

# Interior Point Methods for Optimal Experimental Designs \*

Zhaosong Lu<sup>†</sup>      Ting Kei Pong<sup>‡</sup>

January 25, 2012

## Abstract

In this paper, we propose a primal IP method for solving the optimal experimental design problem with a large class of smooth convex optimality criteria, including A-, D- and  $p$ th mean criterion, and establish its global convergence. We also show that the Newton direction can be computed efficiently when the size of the moment matrix is small relative to the sample size. We compare our IP method with the widely used multiplicative algorithm introduced by Silvey et al. [27]. The computational results show that the IP method generally outperforms the multiplicative algorithm both in speed and solution quality.

**Key words:** Optimal experimental design, A-criterion, c-criterion, D-criterion,  $p$ th mean criterion, interior point methods

## 1 Introduction

In this paper, we consider the optimal experimental design problems on a given finite design space  $\mathcal{X} = \{x_1, \dots, x_n\} \subseteq \mathbb{R}^m$ . In this setting, we consider a coefficient matrix  $K \in \mathbb{R}^{m \times k}$  of full column rank and the moment matrix defined as

$$\mathcal{M}(w) = \sum_{i=1}^n w_i A_i$$

for  $w \in \Omega := \{w : w_i \geq 0, \sum_{i=1}^n w_i = 1\}$ , where  $A_i$  is the expected Fisher information matrix related to  $x_i$ ,  $i = 1, \dots, n$ . As in [37], throughout this paper we assume that  $A_i$ 's are  $m \times m$  real symmetric positive semidefinite matrices and that there exists an  $w \in \Omega$  such that  $\mathcal{M}(w)$  is positive definite. This in particular implies that  $\mathcal{M}(w)$  is positive definite for all positive  $w \in \Omega$ . The optimal experimental design problem can then be formulated as the following minimization problem (see [25, Section 7.10]):

$$\begin{aligned} f^* := \inf_w \quad & \Phi(\mathcal{M}(w)) := \Psi(\mathcal{C}_K(\mathcal{M}(w))) \\ \text{s.t.} \quad & w \in \Omega, \text{ Range}(K) \subseteq \text{Range}(\mathcal{M}(w)), \end{aligned} \tag{1}$$

where  $\Psi$  is a function defined on the set of positive definite matrices and  $\mathcal{C}_K(\mathcal{M}(w))$  is the information matrix defined by  $\mathcal{C}_K(\mathcal{M}(w)) := (K^T(\mathcal{M}(w))^\dagger K)^{-1}$ . Here  $A^\dagger$  denotes the Moore-Penrose pseudoinverse of a matrix  $A$ . The well-definedness of  $\mathcal{C}_K(\mathcal{M}(w))$  is guaranteed by the range inclusion condition in the constraint of (1) and the fact that  $K$  has full column rank [25, Chapter 3]. The function  $\Phi$  in the objective is commonly referred to as an “optimality criterion”. Some classical optimality criteria include (see [25, Chapter 6]):

\*This work was supported in part by an NSERC Discovery Grant.

<sup>†</sup>Department of Mathematics, Simon Fraser University, Burnaby, BC, V5A 1S6, Canada. (email: zhaosong@sfu.ca).

<sup>‡</sup>Department of Combinatorics and Optimization, University of Waterloo, Waterloo, ON, N2L 3G1, Canada. (email: ptingkei@math.uwaterloo.ca).

- (i) A-criterion  $\Phi(X) := \text{tr}(K^T X^\dagger K)$ ;
- (ii) c-criterion  $\Phi(X) := c^T X^\dagger c$ ;
- (iii) D-criterion  $\Phi(X) := \log \det(K^T X^\dagger K)$ ;
- (iv)  $p$ th mean criterion  $\Phi(X) := \text{tr}((K^T X^\dagger K)^{-p})$ .

for some  $p < 0$ ,  $c \in \mathbb{R}^m$  and  $K \in \mathbb{R}^{m \times k}$  of full column rank.

It is easy to observe that c-criterion is just a special case of A-criterion with  $K = c$  and A-criterion is a special case of  $p$ th mean criterion with  $p = -1$ . We shall also mention that  $p$ th mean criterion can be defined more generally to include D-criterion as a special case (see [25, Chapter 6] for details). Furthermore, it can be shown that the constraint set of (1) is convex [25, Section 3.3], and the criteria (i)–(iv) are convex functions in the constraint set (by using [25, Theorem 5.14] and [25, Theorem 6.13], or [24, Proposition IV.14] and [24, Proposition IV.15]). Hence, problem (1) with these criteria is a convex optimization problem. Indeed, it is known that (1) with the above criteria can be reformulated as (possibly nonlinear) semidefinite programming (SDP) problems (see, for example, [14, 7, 9, 23]).

The optimal design problems (1) with the aforementioned criteria usually do not have closed form solutions. Numerous procedures have thus been proposed to solve (1) (see, for example, [13, 36, 4, 5, 35, 6, 18, 8, 24, 2, 32, 12, 33]). Among them, the multiplicative algorithm introduced in [27] has been widely explored. For example, Titterton [28], Pázmán [24], Dette et al. [12] and Harman and Trnovská [19] studied the multiplicative algorithm for D-criterion. In addition, Fellman [15] and Torsney [31] considered the multiplicative algorithm for A-criterion under the assumption that all  $A_i$ 's are rank-one. Recently, Yu [37] studied the multiplicative algorithm for a class of convex optimality criteria and proved its global convergence under some assumptions. Nevertheless, for several commonly used optimality criteria, some of those assumptions may not hold and hence there is no theoretical guarantee for its convergence. Indeed, as observed in [37, Section 5], one of the assumptions does not hold for  $p$ th mean criterion with  $p = -2$ . Moreover, for such a criterion, our numerical experiments in Section 5 demonstrate that the multiplicative algorithm appears not to converge when  $p < -1$ . More details about the multiplicative algorithm for solving (1) are given in Section 2.

In this paper, we consider an alternative approach to solve problem (1). In particular, we study interior point (IP) methods for (1), which are Newton-type methods and can be efficiently applied to a broad class of optimal design problems with moderate-sized matrices  $A_i$ 's. More specifically, we develop a primal IP method for (1) with a large class of convex optimality criteria and establish its global convergence. We also demonstrate how the Newton direction can be computed efficiently when  $n \gg m^2$ , i.e., when the size of  $A_i$ 's is small relative to the sample size. We then compare the IP method with the multiplicative algorithm. The computational results show that the IP method usually outperforms the multiplicative algorithm in both speed and solution quality.

The rest of this paper is organized as follows. In Subsection 1.1, we introduce the notations that are used throughout the paper. In Section 2, we review the multiplicative algorithm and address its convergence. In Sections 3, we propose a primal IP method for solving problem (1) with a large class of convex optimality criteria and address its convergence. In Section 4, we discuss how the IP method can be applied to solve problem (1) with criteria (i)–(iv) and demonstrate how the Newton direction can be computed efficiently when  $n \gg m^2$ . In Section 5, we conduct numerical experiments to test the performance of the method and compare it with the multiplicative algorithm. Finally, we present some concluding remarks in Section 6.

## 1.1 Notations

In this paper, the symbol  $\mathbb{R}_{++}$  denotes the set of all positive real numbers and  $\mathbb{R}^n$  denotes the  $n$ -dimensional Euclidean space. For a vector  $x \in \mathbb{R}^n$  and  $\mathcal{I} \subseteq \{1, \dots, n\}$ ,  $\|x\|$  denotes the Euclidean norm of  $x$ ,  $x_{\mathcal{I}}$  denotes the subvector of  $x$  indexed by  $\mathcal{I}$  and  $\mathcal{D}(x)$  denotes the diagonal matrix whose  $i$ th diagonal entry is  $x_i$  for all  $i$ . For  $\alpha \in \mathbb{R}$  and a vector  $x \in \mathbb{R}^n$  with positive entries,  $x^\alpha$  denotes

the vector whose  $i$ th entry is  $x_i^\alpha$  for all  $i$ . For  $x, y \in \mathbb{R}^n$ ,  $x \circ y$  denotes the Hadamard (entry-wise) product of  $x$  and  $y$ . The letter  $e$  denotes the vector of all ones, whose dimension should be clear from the context. The set of all  $m \times n$  matrices with real entries is denoted by  $\mathbb{R}^{m \times n}$ . For any  $A \in \mathbb{R}^{m \times n}$ ,  $\mathcal{I} \subseteq \{1, \dots, m\}$  and  $\mathcal{J} \subseteq \{1, \dots, n\}$ ,  $a_{ij}$  denotes the  $(i, j)$ th entry of  $A$ ,  $A_{\mathcal{J}}$  denotes the submatrix of  $A$  comprising the columns of  $A$  indexed by  $\mathcal{J}$  and  $A_{\mathcal{I}\mathcal{J}}$  denotes the submatrix of  $A$  comprising the rows and columns of  $A$  indexed by  $\mathcal{I}$  and  $\mathcal{J}$ , respectively. The space of  $n \times n$  symmetric matrices will be denoted by  $\mathcal{S}^n$ . If  $A \in \mathcal{S}^n$  is positive semidefinite (resp., definite), we write  $A \succeq 0$  (resp.,  $A \succ 0$ ). The cone of positive semidefinite (resp., definite) matrices is denoted by  $\mathcal{S}_+^n$  (resp.,  $\mathcal{S}_{++}^n$ ). For  $A, B \in \mathcal{S}^n$ ,  $A \succeq B$  (resp.,  $A \succ B$ ) means  $A - B \succeq 0$  (resp.,  $A - B \succ 0$ ). The trace of a real square matrix  $A$  is denoted by  $\text{tr}(A)$ . We denote by  $I$  the identity matrix, whose dimension should be clear from the context.

A function  $f : \mathcal{S}^n \rightarrow \mathbb{R}$  is said to be increasing (resp., decreasing) if for any  $A \succeq B$ , it holds that

$$f(A) \geq f(B) \quad (\text{resp.}, f(A) \leq f(B)).$$

## 2 The multiplicative algorithm

In this section we review the multiplicative algorithm introduced in [27] for solving problem (1) and discuss its convergence. In particular, we first describe the multiplicative algorithm as follows, which is specified through a power parameter  $\lambda \in (0, 1]$ .

### Multiplicative Algorithm:

1. **Start:** Let a positive  $w^0 \in \Omega$  and  $\lambda \in (0, 1]$  be given.
2. **For**  $k = 0, 1, \dots$

$$w_i^{k+1} = w_i^k \frac{(d_i(w^k))^\lambda}{\sum_{j=1}^n w_j^k (d_j(w^k))^\lambda}, \quad i = 1, \dots, n, \quad (2)$$

where  $d_i(w) = -\text{tr}(\nabla\Phi(\mathcal{M}(w))A_i)$  and  $\nabla\Phi(\mathcal{M}(w))$  is the gradient of  $\Phi$  at  $\mathcal{M}(w)$ .

**End** (for)

**Remark 2.1.** The above algorithm is the same as the one described in [37], in the sense that both algorithms generate exactly the same sequence  $\{w^k\}$  provided the initial points  $w^0$  are identical. ■

We now state a global convergence result recently established by Yu [37, Theorem 2] for the multiplicative algorithm when applied to solve the following problem, which is closely related to (1):

$$\begin{aligned} \text{val} &:= \sup_w -\Phi(\mathcal{M}(w)) \\ \text{s.t.} \quad & w \in \Omega, \mathcal{M}(w) \succ 0. \end{aligned} \quad (3)$$

Observe that (1) and (3) are equivalent (i.e., the optimal value being negative of each other) if there exists an optimal solution  $w^*$  of (1) with  $\mathcal{M}(w^*) \succ 0$ , or if  $\Phi$  is convex in

$$\mathcal{S}_+^m(K) := \{X \in \mathcal{S}_+^m : \text{Range}(K) \subseteq \text{Range}(X)\}$$

and (3) has an optimal solution.

**Proposition 2.1.** Let  $\{w^k\}$  be the sequence generated from the above multiplicative algorithm. Suppose the following assumptions hold:

- (a) for any feasible point  $w$  of (3),  $\nabla\Phi(\mathcal{M}(w)) \preceq 0$  and  $\nabla\Phi(\mathcal{M}(w))A_i \neq 0$  for  $i = 1, \dots, n$ ;
- (b) for any feasible point  $w$  of (3), if  $T(w) \neq w$ , then  $\Phi(\mathcal{M}(T(w))) < \Phi(\mathcal{M}(w))$ , where

$$[T(w)]_i := w_i \frac{(d_i(w))^\lambda}{\sum_{j=1}^n w_j (d_j(w))^\lambda}, \quad i = 1, \dots, n;$$

(c)  $\Phi$  is strictly convex and  $\nabla\Phi$  is continuous in  $\mathcal{S}_{++}^m$ ;

(d) for any  $\{X^k\} \subset \mathcal{S}_{++}^m$ , if  $X^k \rightarrow X^*$  and  $\{\Phi(X^k)\}$  is decreasing, then  $X^* \succ 0$ .

Then  $\Phi(\mathcal{M}(w^k)) \rightarrow -\text{val}$  monotonically, and moreover, any accumulation point of  $\{w^k\}$  is an optimal solution of (3).

**Remark 2.2.** Notice that the assumptions in the above proposition imply that any accumulation point  $w^*$  of  $\{w^k\}$  satisfies  $\mathcal{M}(w^*) \succ 0$ . Hence, if the assumptions in Proposition 2.1 hold and  $\Phi$  is convex in  $\mathcal{S}_+^m(K)$ , then (1) is equivalent to (3) and any accumulation point of the sequence  $\{w^k\}$  generated from the above multiplicative algorithm solves (1). ■

Using Proposition 2.1 and some technical results developed in [37], one can establish the convergence of the above multiplicative algorithm when applied to problem (1) with A-, D- and pth mean criterion for  $p \in (-1, 0)$  and  $K = I$ , which is summarized as follows.

**Corollary 2.1.** Assume that  $K = I$  and  $A_i \neq 0$  for  $i = 1, \dots, n$ . Then the multiplicative algorithm converges for any  $\lambda \in (0, 1]$  when applied to problem (1) with D- and pth mean criterion for  $p \in (-1, 0)$ . Also, it converges for A-criterion when  $\lambda \in (0, 1)$ .

As seen from Proposition 2.1 and Corollary 2.1, the multiplicative algorithm converges for a large class of optimality criteria  $\Phi$ . Nevertheless, for some important convex optimality criteria, the assumptions stated in Proposition 2.1 may not hold and hence there is no theoretical guarantee for its convergence. Indeed, as observed in [37, Section 5], the assumption (b) with  $\lambda = 1$  does not hold for pth mean criterion with  $p = -2$ . Moreover, for such a criterion, our numerical experiments in Section 5 demonstrate that the multiplicative algorithm appears not to converge when  $p < -1$ .

Due to the aforementioned potential drawbacks of the multiplicative algorithm, we will propose an IP method for solving problem (1) with a broad class of optimality criteria  $\Phi$  including A-, D- and pth mean criterion in subsequent sections.

### 3 IP method for a class of convex optimality criteria

In this section, we propose a primal IP method for solving (1) with a class of convex optimality criteria  $\Phi = \Psi \circ \mathcal{C}_K$ . We make the following assumption on  $\Psi$  throughout this section.

**Assumption 3.1.** The function  $\Psi$  is convex, decreasing, twice continuously differentiable and bounded below in  $\mathcal{S}_{++}^m$ . Moreover, for any bounded sequences  $\{X^k\} \subseteq \mathcal{S}_{++}^m$  with  $\lambda_{\min}(X^k) \rightarrow 0$ , one has  $\Psi(X^k) \rightarrow \infty$ .

**Remark 3.1.** We now make some brief comments on the above assumptions.

- (a) Assumption 3.1 is fairly reasonable. Indeed, all optimality criteria described in Section 1 satisfy this assumption.
- (b) Since the feasible set is not necessarily closed, problem (1) with a general convex optimality criterion may not have an optimal solution. However, when the optimality criterion satisfies Assumption 3.1, it must have an optimal solution as shown in Theorem 3.1(a). We refer the readers to [25, Chapter 5] for more discussion on conditions guaranteeing existence of solutions for problem (1).
- (c) In contrast to Proposition 2.1, we do not require the existence of a positive definite optimal moment matrix  $\mathcal{M}(w^*)$ . Indeed, Assumption 3.1 may hold even when problem (1) does not have a positive definite optimal moment matrix. For instance, the design problem

$$\begin{aligned} \min_{w, X} \quad & \begin{pmatrix} 1 \\ 0 \end{pmatrix}^T X^\dagger \begin{pmatrix} 1 \\ 0 \end{pmatrix} \\ \text{s.t.} \quad & X = \begin{pmatrix} w_1 & 0 \\ 0 & w_2 \end{pmatrix}, w_1 + w_2 = 1, w_1, w_2 \geq 0, \\ & \begin{pmatrix} 1 \\ 0 \end{pmatrix} \in \text{Range}(X), \end{aligned}$$

has a unique optimal solution at  $(w_1, w_2) = (1, 0)$ . The corresponding optimal moment matrix is not positive definite; thus, the assumption (d) of Proposition 2.1 does not hold. However, it is easy to check that Assumption 3.1 is satisfied for this design problem (with  $\Psi(t) = 1/t$ ). In general, the assumption (d) of Proposition 2.1 does not hold when  $K$  is not invertible, while our Assumption 3.1 is independent of  $K$ . ■

Under Assumption 3.1, it is not hard to show that the function  $\Phi(\mathcal{M}(\cdot))$  is bounded below on the feasible set of (1). Also, it is routine to show that the function  $\Phi$  is twice continuously differentiable in  $\mathcal{S}_{++}^m$ . Furthermore, it can be shown that  $\Phi$  is convex in  $\mathcal{S}_+^m(K)$  by considering suitable Schur complements (see, for example, [23, Section 6]). We include a short proof below for the convenience of readers. Before proceeding, we state the following well-known fact, which concerns the Schur complement of a positive semidefinite submatrix (see, for example, [25, Lemma 3.12]).

**Lemma 3.1.** *Let  $A \in \mathcal{S}^k$ ,  $B \in \mathbb{R}^{m \times k}$  and  $C \in \mathcal{S}^m$ . Then the matrix  $\begin{pmatrix} A & B^T \\ B & C \end{pmatrix}$  is positive semidefinite if and only if  $A \succeq B^T C^\dagger B$ ,  $C \succeq 0$  and  $\text{Range}(B) \subseteq \text{Range}(C)$ .*

**Proposition 3.1.** *The optimality criterion  $\Phi$  is convex in  $\mathcal{S}_+^m(K)$ .*

*Proof.* First of all, it can be shown that the set  $\mathcal{S}_+^m(K)$  is convex (see, for example, [25, Section 3.3]). In addition, notice that for any  $X \in \mathcal{S}_+^m(K)$ , we have

$$\begin{aligned} \Phi(X) &= \Psi((K^T X^\dagger K)^{-1}) = \inf_U \{ \Psi(U) : (K^T X^\dagger K)^{-1} \succeq U \succ 0 \} \\ &= \inf_U \{ \Psi(U) : U^{-1} \succeq K^T X^\dagger K, U \succ 0 \} \\ &= \inf_U \left\{ \Psi(U) : \begin{pmatrix} U^{-1} & K^T \\ K & X \end{pmatrix} \succeq 0, U \succ 0 \right\} \\ &= \inf_U \{ \Psi(U) : X \succeq K U K^T, U \succ 0 \}, \end{aligned} \tag{4}$$

where the second equality follows from the fact that  $\Psi$  is decreasing, the fourth and last equalities follow from Lemma 3.1, while the third equality holds because  $K^T X^\dagger K$  is invertible for  $X \in \mathcal{S}_+^m(K)$  when  $K$  has full column rank. Convexity of  $\Phi$  in  $\mathcal{S}_+^m(K)$  now follows from [26, Theorem 5.7]. ■

Observe that  $\mathcal{M}(w) \succ 0$  whenever  $w > 0$ . Thus, under Assumption 3.1, the function  $\Phi$  is twice continuously differentiable for any positive  $w \in \Omega$ . It is hence natural to develop an IP method to solve (1) since such a method keeps all iterates in the relative interior of  $\Omega$  until convergence. To proceed, we first reformulate the problem by eliminating the equality constraint. The resulting equivalent problem is given by

$$\begin{aligned} f^* &= \inf_{\tilde{w}} f(\tilde{w}) := \Phi(\mathcal{M}(P\tilde{w} + q)) \\ \text{s.t. } & e^T \tilde{w} \leq 1, \tilde{w} \geq 0, \\ & \text{Range}(K) \subseteq \text{Range}(\mathcal{M}(P\tilde{w} + q)), \end{aligned} \tag{5}$$

where  $P \in \mathbb{R}^{n \times (n-1)}$  and  $q \in \mathbb{R}^n$  are such that

$$P\tilde{w} + q = \begin{pmatrix} \tilde{w} \\ 1 - e^T \tilde{w} \end{pmatrix} \quad \forall \tilde{w} \in \mathbb{R}^{n-1}. \tag{6}$$

We next develop an IP method for solving problem (5) instead. First, we need to build a suitable barrier function. Given any  $\tilde{w} > 0$  satisfying  $e^T \tilde{w} < 1$ , one can observe that  $P\tilde{w} + q \succ 0$  and hence  $\mathcal{M}(P\tilde{w} + q) \succ 0$ , which leads to  $\text{Range}(K) \subseteq \text{Range}(\mathcal{M}(P\tilde{w} + q))$ . This implies that any barrier function that takes into account the first two inequality constraints of (5) is sufficient

for the development of IP method. Here we naturally choose the logarithmic barrier function and then solve the barrier subproblem in the form of

$$\min_{\tilde{w}} f_{\mu}(\tilde{w}) := f(\tilde{w}) - \mu \sum_{i=1}^{n-1} \log(\tilde{w}_i) - \mu \log(1 - e^T \tilde{w}) \quad (7)$$

for a sequence of parameters  $\mu \downarrow 0$ . In view of Assumption 3.1, we see that any level set of  $f_{\mu}$  is compact. Moreover,  $f_{\mu}$  is strictly convex. Thus, there exists a unique minimizer to (7) for any  $\mu > 0$ . Furthermore, it follows from Assumption 3.1 that  $f_{\mu}$  is twice continuously differentiable and its Hessian is positive definite in its domain. Therefore, problem (7) can be suitably solved by the Newton's method with a line search whose stepsize is chosen by Armijo rule.

We are now ready to present our IP method for solving problem (5).

### Primal IP Method:

1. **Start:** Let a strictly feasible  $\tilde{w}^0$ ,  $0 < \beta, \gamma, \eta, \sigma < 1$  and  $\mu_1 > 0$  be given. Let  $\epsilon(\mu)$  be an increasing function of  $\mu$  so that  $\lim_{\mu \downarrow 0} \epsilon(\mu) = 0$ . Set  $\tilde{w} = \tilde{w}^0$  and  $k = 1$ .

2. **While**  $\|\nabla f_{\mu_k}(\tilde{w})\| > \epsilon(\mu_k)$  **do**

(a) Compute the Newton direction

$$d := -(\nabla^2 f_{\mu_k}(\tilde{w}))^{-1} \nabla f_{\mu_k}(\tilde{w}). \quad (8)$$

(b) Let  $\alpha_{\max}(\tilde{w}) := \max\{\alpha : \tilde{w}[\alpha] \geq 0, e^T \tilde{w}[\alpha] \leq 1\}$ , where  $\tilde{w}[\alpha] := \tilde{w} + \alpha d$ .

(c) Let  $\alpha$  be the largest element of  $\{\bar{\alpha}(\tilde{w}), \beta \bar{\alpha}(\tilde{w}), \beta^2 \bar{\alpha}(\tilde{w}), \dots\}$  satisfying

$$f_{\mu_k}(\tilde{w}[\alpha]) \leq f_{\mu_k}(\tilde{w}) + \sigma \alpha (\nabla f_{\mu_k}(\tilde{w}))^T d,$$

where  $\bar{\alpha}(\tilde{w}) := \min\{1, \eta \alpha_{\max}(\tilde{w})\}$ .

(d) Set  $\tilde{w} \leftarrow \tilde{w}[\alpha]$ .

**End** (while)

3. Set  $\tilde{w}^k \leftarrow \tilde{w}$ ,  $\mu_{k+1} \leftarrow \gamma \mu_k$ ,  $k \leftarrow k + 1$ , and go to step 2.

In standard convergence analysis of IP methods, the feasible sets are usually assumed to be closed and the objective functions are twice continuously differentiable in a neighborhood of the feasible sets (see, for example, [17]). Nevertheless, these two conditions do not necessarily hold for our problem (5). In particular, the objective function is not necessarily continuous up to the boundary of the feasible region [25, Section 3.16]. Hence, it is not immediately clear the sequence generated by our method will accumulate at a global minimizer of (5). Thus, we analyze convergence of our IP method below. We first present convergence analysis regarding the outer iterations of our IP method and leave the discussion on the convergence of its inner iterations to the end of this section.

For notational convenience, in the remainder of this section, we associate with each  $\tilde{w} \in \mathbb{R}^{n-1}$  a unique  $w \in \mathbb{R}^n$  by letting  $w := P\tilde{w} + q$ . Analogously, we associate with each  $w \in \mathbb{R}^n$  a unique  $\tilde{w} \in \mathbb{R}^{n-1}$  by letting  $\tilde{w}_i = w_i$  for  $i = 1, \dots, n-1$ . Also, we let  $\Phi_{\mathcal{M}}(w) := \Phi(\mathcal{M}(w))$ .

We first observe that if problem (1) has an optimal solution  $w^*$  with  $\mathcal{M}(w^*) \succ 0$ , then there exists a Lagrange multiplier  $u^* \geq 0$  such that  $(w^*, u^*)$  satisfies the following KKT system:

$$\begin{aligned} P^T (\nabla \Phi_{\mathcal{M}}(w) - u) &= 0, \\ e^T w &= 1, \\ u \circ w &= 0, \\ (w, u) &\geq 0. \end{aligned} \quad (9)$$

Given a strictly feasible point  $\tilde{w} \in \mathbb{R}^{n-1}$  of problem (7), we notice that

$$\nabla f_\mu(\tilde{w}) = P^T(\nabla \Phi_{\mathcal{M}}(w) - \mu w^{-1}). \quad (10)$$

Then it is not hard to observe that for each  $\mu > 0$ , the  $w$  associated with the approximate solution  $\tilde{w}$  of (7) obtained by the Newton's method detailed in step 2 above together with  $u := \mu w^{-1}$  satisfies the following perturbed KKT system:

$$\begin{aligned} P^T(\nabla \Phi_{\mathcal{M}}(w) - u) &= v, \\ e^T w &= 1, \\ u \circ w &= \mu e, \\ (w, u) &> 0 \end{aligned} \quad (11)$$

for some  $v \in \mathbb{R}^{n-1}$ . In order to analyze the convergence of our IP method, we will study the limiting behavior of the solutions of system (11) as  $(\mu, v) \rightarrow (0_+, 0)$ , that is,  $(\mu, v) \rightarrow (0, 0)$  with  $\mu > 0$ .

We first claim that system (11) has a unique solution for any  $(\mu, v) \in \mathbb{R}_{++} \times \mathbb{R}^{n-1}$ . Indeed, it is easy to observe that  $(w, u)$  is a solution of (11) if and only if  $\tilde{w} \in \mathbb{R}^{n-1}$  is an optimal solution of

$$\min_{\tilde{w}} f_\mu(\tilde{w}) - v^T \tilde{w}. \quad (12)$$

Since the objective function of (12) is strictly convex and it has compact level sets, problem (12) has a unique optimal solution, which immediately implies that system (11) has a unique solution. From now on, we denote by  $(w(\mu, v), u(\mu, v))$  the unique solution of (11) and by  $\tilde{w}(\mu, v)$  the vector obtained from  $w(\mu, v)$  by dropping the last entry for all  $(\mu, v) \in \mathbb{R}_{++} \times \mathbb{R}^{n-1}$ . It is clear that  $\tilde{w}(\mu, v)$  is the unique optimal solution of (12). We next establish the limiting behavior of  $w(\mu, v)$  as  $(\mu, v) \rightarrow (0_+, 0)$ .

**Theorem 3.1.** *Let  $(w(\mu, v), u(\mu, v))$  be defined above for  $(\mu, v) \in \mathbb{R}_{++} \times \mathbb{R}^{n-1}$ . Then the following statements hold:*

- (a)  $\lim_{(\mu, v) \rightarrow (0_+, 0)} \Phi(\mathcal{M}(w(\mu, v))) = f^*$  and any accumulation point of  $w(\mu, v)$  as  $(\mu, v) \rightarrow (0_+, 0)$  is an optimal solution of (1).
- (b) Suppose that problem (1) has an optimal solution  $w^*$  with  $\mathcal{M}(w^*) \succ 0$ . Then any accumulation point of  $w(\mu, v)$  as  $(\mu, v) \xrightarrow{\Xi_C} (0, 0)$ , i.e.,  $(\mu, v) \rightarrow (0, 0)$  with  $(\mu, v) \in \Xi_C := \{(\mu, v) : \|v\|_\infty < C\mu\}$  for some given  $C > 0$ , is an optimal solution of (1) with maximum cardinality.
- (c) Suppose that problem (1) has an optimal solution  $w^*$  with  $\mathcal{M}(w^*) \succ 0$ . Let  $u^*$  be the associated Lagrange multiplier satisfying (9). Assume that  $|\mathcal{B}| + |\mathcal{N}| = n$ , where  $\mathcal{B} := \{i : w_i^* > 0\}$  and  $\mathcal{N} := \{i : u_i^* > 0\}$ . Suppose further that  $[\nabla^2 \Phi_{\mathcal{M}}(w)]_{\mathcal{B}\mathcal{B}} \succ 0$  for any  $w \in \Omega$  satisfying  $w_{\mathcal{B}} > 0$  and  $w_{\mathcal{N}} = 0$ . Then  $w(\mu, v)$  converges to an optimal solution of (1) with maximum cardinality as  $(\mu, v) \rightarrow (0_+, 0)$ .

*Proof.* We first prove part (a). Let

$$\bar{f}^* := \inf_w \{\Phi_{\mathcal{M}}(w) : e^T w = 1, w > 0\}. \quad (13)$$

We first show that  $\lim_{(\mu, v) \rightarrow (0_+, 0)} \Phi_{\mathcal{M}}(w(\mu, v)) = \bar{f}^*$ .

Given an arbitrary  $\epsilon > 0$ , there exists a positive  $\tilde{w}$  satisfying  $e^T \tilde{w} < 1$  such that  $f(\tilde{w}) < \bar{f}^* + \epsilon/2$ . Recall that  $\tilde{w}(\mu, v)$  is the unique optimal solution of (12). Then we have that for any  $v \in \mathbb{R}^{n-1}$ ,

$$f_\mu(\tilde{w}(\mu, v)) - v^T \tilde{w}(\mu, v) \leq f_\mu(\tilde{w}) - v^T \tilde{w}. \quad (14)$$

On the other hand, note that  $\tilde{w}(\mu, v) > 0$  and  $e^T \tilde{w}(\mu, v) < 1$ . Hence,

$$-\sum_{i=1}^{n-1} \log(\tilde{w}_i(\mu, v)) - \log(1 - e^T \tilde{w}(\mu, v)) > 0$$

and  $f(\tilde{w}(\mu, v)) \geq \bar{f}^*$ . In view of these inequalities, (14) and the fact that  $\|\tilde{w}(\mu, v)\|_1 \leq 1$  and  $\|\tilde{w}\|_1 \leq 1$ , one can obtain that for any  $(\mu, v) \in \mathfrak{R}_{++} \times \mathfrak{R}^{n-1}$ ,

$$\begin{aligned} \bar{f}^* \leq f(\tilde{w}(\mu, v)) &= f_\mu(\tilde{w}(\mu, v)) + \mu \sum_{i=1}^{n-1} \log(\tilde{w}_i(\mu, v)) + \mu \log(1 - e^T \tilde{w}(\mu, v)) \\ &\leq f_\mu(\tilde{w}(\mu, v)) \leq f_\mu(\tilde{w}) + v^T \tilde{w}(\mu, v) - v^T \tilde{w} \\ &\leq f(\tilde{w}) - \mu \sum_{i=1}^{n-1} \log(\tilde{w}_i) - \mu \log(1 - e^T \tilde{w}) + 2\|v\|_\infty \\ &\leq \bar{f}^* + \frac{\epsilon}{2} - \mu \sum_{i=1}^{n-1} \log(\tilde{w}_i) - \mu \log(1 - e^T \tilde{w}) + 2\|v\|_\infty. \end{aligned}$$

Thus, there exists some  $\delta > 0$  such that  $\bar{f}^* \leq f(\tilde{w}(\mu, v)) \leq \bar{f}^* + \epsilon$  whenever  $\|(\mu, v)\| < \delta$ ,  $\mu > 0$ . Hence,  $\Phi_{\mathcal{M}}(w(\mu, v)) = f(\tilde{w}(\mu, v)) \rightarrow \bar{f}^*$  as  $(\mu, v) \rightarrow (0_+, 0)$ .

We next show that  $f^* = \bar{f}^*$ . Clearly,  $f^* \leq \bar{f}^*$ . We now suppose for contradiction that  $f^* < \bar{f}^*$ . By the definitions of  $f^*$  and  $\bar{f}^*$ , there exist  $w^1$  and  $w^2$  which are feasible points of (1) and (13), respectively, so that  $\Phi_{\mathcal{M}}(w^1) < (f^* + \bar{f}^*)/2$  and  $\Phi_{\mathcal{M}}(w^2) < \bar{f}^* + (\bar{f}^* - f^*)/2$ . Let  $w = (w^1 + w^2)/2$ . Clearly,  $w > 0$ ,  $e^T w = 1$  and  $\text{Range}(K) \subseteq \text{Range}(\mathcal{M}(w))$  due to  $\mathcal{M}(w) \succ 0$ . By convexity of  $\Phi$  in  $\mathcal{S}_+^n(K)$ , we obtain that  $\Phi_{\mathcal{M}}(w) \leq (\Phi_{\mathcal{M}}(w^1) + \Phi_{\mathcal{M}}(w^2))/2 < \bar{f}^*$ , which is a contradiction to the definition of  $\bar{f}^*$ . Thus,  $\lim_{(\mu, v) \rightarrow (0_+, 0)} \Phi_{\mathcal{M}}(w(\mu, v)) = \bar{f}^* = f^*$ .

Now suppose that  $w^*$  is an accumulation point of  $w(\mu, v)$  as  $(\mu, v) \rightarrow (0_+, 0)$ . We next show that  $w^*$  is an optimal solution of (1). Indeed, it follows from (4) that for any feasible point  $w$  of (1),

$$\Phi_{\mathcal{M}}(w) = \inf_U \{ \Psi(U) : \mathcal{M}(w) \succeq KU K^T, U \succ 0 \}. \quad (15)$$

In view of (15), for each  $(\mu, v) \in \mathfrak{R}_{++} \times \mathfrak{R}^{n-1}$ , there exists  $U(\mu, v) \succ 0$  such that

$$\Phi_{\mathcal{M}}(w(\mu, v)) + \|(\mu, v)\| > \Psi(U(\mu, v)) \quad \text{and} \quad \mathcal{M}(w(\mu, v)) \succeq KU(\mu, v)K^T. \quad (16)$$

From the second relation in (16), we see that  $\text{tr}(\mathcal{M}(w(\mu, v))) \geq \lambda_{\min}(K^T K) \text{tr}(U(\mu, v))$ , from which it follows that  $U(\mu, v)$  is bounded and thus it has an accumulation point as  $(\mu, v) \rightarrow (0_+, 0)$ . Let  $U^*$  be such an accumulation point. In view of the first relation in (16) and the assumption on  $\Psi$ , we see that  $U^* \succ 0$ . Moreover, we obtain by taking limit in (16) upon  $(\mu, v) \rightarrow (0_+, 0)$  that

$$\lim_{(\mu, v) \rightarrow (0_+, 0)} \Phi_{\mathcal{M}}(w(\mu, v)) \geq \Psi(U^*), \quad \mathcal{M}(w^*) \succeq KU^* K^T. \quad (17)$$

The second relation in (17) together with Lemma 3.1 implies that

$$\mathcal{M}(w^*) \succeq KU^* K^T \Rightarrow \begin{pmatrix} U^{*-1} & K^T \\ K & \mathcal{M}(w^*) \end{pmatrix} \succeq 0 \Rightarrow \text{Range}(K) \subseteq \text{Range}(\mathcal{M}(w^*)).$$

Hence,  $w^*$  is a feasible point of (1). In view of (15), the first relation in (17) and the result  $\lim_{(\mu, v) \rightarrow (0_+, 0)} \Phi_{\mathcal{M}}(w(\mu, v)) = f^*$ , we have

$$\Phi_{\mathcal{M}}(w^*) \leq \Psi(U^*) \leq \lim_{(\mu, v) \rightarrow (0_+, 0)} \Phi_{\mathcal{M}}(w(\mu, v)) = f^*.$$

Thus,  $w^*$  is an optimal solution of (1). This proves part (a).



We now show that part (b) holds. Let  $w^*$  be an optimal solution of (1) with maximum cardinality. Then it follows immediately from assumption that  $\mathcal{M}(w^*) \succ 0$ . Thus, there exists a corresponding Lagrange multiplier  $u^*$  so that  $(w^*, u^*)$  satisfies (9). Let  $\tilde{w}^*$  be the vector obtained from  $w^*$  by dropping the last entry. In view of (6) and the first equation of (9) and (11), we observe that for any  $(\mu, v) \in \Xi_C$ ,

$$\begin{aligned} & (w(\mu, v) - w^*)^T (u(\mu, v) - u^*) \\ &= (P\tilde{w}(\mu, v) - P\tilde{w}^*)^T (u(\mu, v) - u^*) \\ &= (\tilde{w}(\mu, v) - \tilde{w}^*)^T P^T (\nabla \Phi_{\mathcal{M}}(w(\mu, v)) - \nabla \Phi_{\mathcal{M}}(w^*)) - (\tilde{w}(\mu, v) - \tilde{w}^*)^T v \\ &= (w(\mu, v) - w^*)^T (\nabla \Phi_{\mathcal{M}}(w(\mu, v)) - \nabla \Phi_{\mathcal{M}}(w^*)) - (\tilde{w}(\mu, v) - \tilde{w}^*)^T v \\ &\geq -2C\mu, \end{aligned}$$

where the last inequality holds since  $\Phi$  is convex in  $\mathcal{S}_{++}^m$ ,  $w(\mu, v), w^* \in \Omega$  and  $\|v\|_{\infty} < C\mu$ . Using this inequality and the third equation in (9) and (11), we see that

$$w^{*T} u(\mu, v) + w(\mu, v)^T u^* \leq w^{*T} u^* + w(\mu, v)^T u(\mu, v) + 2C\mu = (2C + n)\mu. \quad (18)$$

Dividing both sides of the above inequality by  $\mu$  and using the third equation of (11), we obtain that

$$\sum_{i=1}^n \frac{w_i^*}{w_i(\mu, v)} + \sum_{i=1}^n \frac{u_i^*}{u_i(\mu, v)} \leq 2C + n. \quad (19)$$

Since  $(w^*, u^*) \geq 0$  and  $(w(\mu, v), u(\mu, v)) > 0$ , it follows from (19) that for all  $i$ ,

$$w_i(\mu, v) \geq \frac{w_i^*}{2C + n}, \quad u_i(\mu, v) \geq \frac{u_i^*}{2C + n}. \quad (20)$$

It immediately implies that the  $i$ th entry of any accumulation point  $w^\diamond$  of  $w(\mu, v)$  as  $(\mu, v) \xrightarrow{\Xi_C} (0, 0)$  must be positive whenever  $w_i^* > 0$ . Since  $w^\diamond$  is an optimal solution of (1) by part (a), we conclude that part (b) holds.

Finally, we show that part (c) holds. By assumption,  $(w^*, u^*)$  is a solution of (9) with  $|\mathcal{B}| + |\mathcal{N}| = n$ . Notice from the third equation of (9) that  $\mathcal{B} \cap \mathcal{N} = \emptyset$ . Thus,  $\mathcal{B}$  and  $\mathcal{N}$  form a partition for  $\{1, \dots, n\}$ . We first show that when  $(\mu, v) \in \Xi_C$  and  $\mu$  is sufficiently small,

$$\begin{aligned} w_{\mathcal{N}}(\mu, v) &= O(\mu), \quad w_i(\mu, v) = \Theta(1), \quad i \in \mathcal{B}, \\ u_{\mathcal{B}}(\mu, v) &= O(\mu), \quad u_i(\mu, v) = \Theta(1), \quad i \in \mathcal{N}. \end{aligned} \quad (21)$$

Since  $w_{\mathcal{B}}^* > 0$  and  $u_{\mathcal{N}}^* > 0$ , it follows from (18) that  $w_{\mathcal{N}}(\mu, v) = O(\mu)$  and  $u_{\mathcal{B}}(\mu, v) = O(\mu)$  for all  $(\mu, v) \in \Xi_C$ . In addition, in view of (20) and the fact that  $w(\mu, v) \in \Omega$ , one can immediately see that  $w_i(\mu, v) = \Theta(1)$  for all  $i \in \mathcal{B}$  and  $(\mu, v) \in \Xi_C$ . We now show that when  $(\mu, v) \in \Xi_C$  and  $\mu$  is sufficiently small,  $u_i(\mu, v) = \Theta(1)$  for all  $i \in \mathcal{N}$ . Indeed, using the first equation of (11), we obtain that

$$P^T \nabla \Phi_{\mathcal{M}}(w(\mu, v)) - (P^T)_{\mathcal{B}} u_{\mathcal{B}}(\mu, v) - (P^T)_{\mathcal{N}} u_{\mathcal{N}}(\mu, v) = v.$$

Since  $w^* \in \Omega$ , we know that  $|\mathcal{B}| \geq 1$ , which together with the definition of  $P$  implies that  $(P^T)_{\mathcal{N}}$  has full column rank. It then follows from the above equation that

$$u_{\mathcal{N}}(\mu, v) = ([ (P^T)_{\mathcal{N}} ]^T [ (P^T)_{\mathcal{N}} ]^{-1} [ (P^T)_{\mathcal{N}} ]^T (P^T \nabla \Phi_{\mathcal{M}}(w(\mu, v)) - (P^T)_{\mathcal{B}} u_{\mathcal{B}}(\mu, v) - v)). \quad (22)$$

Recall from above that  $w_i(\mu, v) = \Theta(1)$  for all  $i \in \mathcal{B}$  and  $(\mu, v) \in \Xi_C$ . Using this result and the fact that  $w(\mu, v) \in \Omega$  and  $\mathcal{M}(w^*) \succ 0$ , it is not hard to see that  $\{\mathcal{M}(w(\mu, v)) : (\mu, v) \in \Xi_C\}$  is included in a compact set in  $\mathcal{S}_{++}^m$ . Hence,  $P^T \nabla \Phi_{\mathcal{M}}(w(\mu, v))$  is bounded for all  $(\mu, v) \in \Xi_C$ .

Further, in view of (20), (22) and the result that  $u_{\mathcal{B}}(\mu, v) = O(\mu)$  for all  $(\mu, v) \in \Xi_C$ , we easily see that  $u_i(\mu, v) = \Theta(1)$  for all  $i \in \mathcal{N}$  when  $(\mu, v) \in \Xi_C$  and  $\mu$  is sufficiently small.

For all  $(\mu, v) \in \mathfrak{R}_{++} \times \mathfrak{R}^{n-1}$ , let

$$w^+(\mu, v) := \left( w_{\mathcal{B}}(\mu, v), \frac{1}{\mu} w_{\mathcal{N}}(\mu, v) \right), \quad u^+(\mu, v) := \left( \frac{1}{\mu} u_{\mathcal{B}}(\mu, v), u_{\mathcal{N}}(\mu, v) \right), \quad (23)$$

$$I_1(\mu) = \begin{pmatrix} I & 0 \\ 0 & \mu I \end{pmatrix}, \quad I_2(\mu) = \begin{pmatrix} \mu I & 0 \\ 0 & I \end{pmatrix}.$$

Then it follows from (11) that  $(w^+(\mu, v), u^+(\mu, v))$  is the unique solution of

$$F(w, u, \mu, v) := \begin{pmatrix} P^T(\nabla \Phi_{\mathcal{M}}(I_1(\mu)w) - I_2(\mu)u) - v \\ e^T(I_1(\mu)w) - 1 \\ u \circ w - e \end{pmatrix} = 0, \quad (w, u) > 0 \quad (24)$$

for all  $(\mu, v) \in \mathfrak{R}_{++} \times \mathfrak{R}^{n-1}$ . In view of (21), we know that when  $(\mu, v) \in \Xi_C$  and  $\mu$  is sufficiently small,

$$w_{\mathcal{N}}^+(\mu, v) = O(1), \quad w_i^+(\mu, v) = \Theta(1), \quad i \in \mathcal{B},$$

$$u_{\mathcal{B}}^+(\mu, v) = O(1), \quad u_i^+(\mu, v) = \Theta(1), \quad i \in \mathcal{N}.$$

Thus,  $(w^+(\mu, v), u^+(\mu, v))$  has accumulation points as  $(\mu, v) \xrightarrow{\Xi_C} (0, 0)$ . Let  $(w^\diamond, u^\diamond)$  be one such accumulation point. Clearly,  $(w^\diamond, u^\diamond) > 0$  due to the third equation and the inequality in (24). We will show below that  $(w^+(\mu, v), u^+(\mu, v)) \rightarrow (w^\diamond, u^\diamond)$  as  $(\mu, v) \rightarrow (0_+, 0)$ .

First, it is easy to see that  $(w^\diamond, u^\diamond, 0, 0)$  satisfies (24). Since  $\mathcal{M}(w^*) \succ 0$  and  $w^\diamond > 0$ , it is not hard to verify  $\mathcal{M}(I_1(0)w^\diamond) \succ 0$ . Using this result and Assumption 3.1, we observe that  $F$  is continuously differentiable in a neighborhood of  $(w^\diamond, u^\diamond, 0, 0)$ . The Jacobian matrix of  $F$  with respect to  $(w, u)$  at  $(w^\diamond, u^\diamond, 0, 0)$  is given by

$$J := \begin{pmatrix} P^T \nabla^2 \Phi_{\mathcal{M}}(I_1(0)w^\diamond) I_1(0) & -P^T I_2(0) \\ e^T I_1(0) & 0 \\ \mathcal{D}(u^\diamond) & \mathcal{D}(w^\diamond) \end{pmatrix}.$$

We next show that the Jacobian matrix  $J$  is nonsingular. It suffices to show that the linear system  $J \begin{pmatrix} \Delta w \\ \Delta u \end{pmatrix} = 0$ , or equivalently,

$$\begin{aligned} P^T \nabla^2 \Phi_{\mathcal{M}}(I_1(0)w^\diamond) I_1(0) \Delta w - P^T I_2(0) \Delta u &= 0, \\ e^T(I_1(0) \Delta w) &= 0, \\ u^\diamond \circ \Delta w + w^\diamond \circ \Delta u &= 0, \end{aligned} \quad (25)$$

has only zero solution. Indeed, we observe that the null space of  $P^T$  is spanned by  $e$ . It then follows from the first equation of (25) that

$$\nabla^2 \Phi_{\mathcal{M}}(I_1(0)w^\diamond) I_1(0) \Delta w - I_2(0) \Delta u = \lambda e, \quad (26)$$

for some  $\lambda \in \mathfrak{R}$ . Multiplying both sides of (26) by  $(I_1(0) \Delta w)^T$  and making use of the second equation of (25) and the fact  $I_1(0) I_2(0) = 0$ , we arrive at

$$\Delta w^T I_1(0) \nabla^2 \Phi_{\mathcal{M}}(I_1(0)w^\diamond) I_1(0) \Delta w = 0,$$

which is equivalent to

$$\Delta w_{\mathcal{B}}^T [\nabla^2 \Phi_{\mathcal{M}}(I_1(0)w^\diamond)]_{\mathcal{B}\mathcal{B}} \Delta w_{\mathcal{B}} = 0.$$

This together with the assumption implies that  $\Delta w_{\mathcal{B}} = 0$ . Using this result, the third equation of (25) and the fact  $w^\diamond > 0$ , we see that  $\Delta u_{\mathcal{B}} = 0$ . Substituting  $\Delta w_{\mathcal{B}} = 0$  into (26), we obtain that

$-I_2(0)\Delta u = \lambda e$ , which together with the definition of  $I_2(0)$  and  $|\mathcal{B}| \geq 1$  implies  $\lambda = 0$  and hence  $\Delta u_{\mathcal{N}} = 0$ . Using this result, the third equation of (25) and the fact  $u^\diamond > 0$ , we see that  $\Delta w_{\mathcal{N}} = 0$ . Thus, we have shown that  $\Delta w = \Delta u = 0$ . Hence, the Jacobian matrix  $J$  is nonsingular.

By applying the implicit function theorem to (24), we conclude that there exists  $\epsilon > 0$ , a neighborhood  $\mathcal{U}$  of 0 and a continuously differentiable function  $(w^\dagger(\mu, v), u^\dagger(\mu, v))$  defined on  $(-\epsilon, \epsilon) \times \mathcal{U}$  such that

$$\begin{aligned} F(w^\dagger(\mu, v), u^\dagger(\mu, v), \mu, v) &= 0, \quad (w^\dagger(\mu, v), u^\dagger(\mu, v)) > 0 \quad \forall (\mu, v) \in (-\epsilon, \epsilon) \times \mathcal{U}, \\ \lim_{(\mu, v) \rightarrow 0} (w^\dagger(\mu, v), u^\dagger(\mu, v)) &= (w^\diamond, u^\diamond). \end{aligned}$$

Recall that system (24) has a unique solution  $(w^+(\mu, v), u^+(\mu, v))$  for all  $(\mu, v) \in \mathfrak{R}_{++} \times \mathfrak{R}^{n-1}$ . Therefore,  $(w^+(\mu, v), u^+(\mu, v)) = (w^\dagger(\mu, v), u^\dagger(\mu, v))$  for all  $(\mu, v) \in (0, \epsilon) \times \mathcal{U}$ . It then follows that

$$\lim_{(\mu, v) \rightarrow (0_+, 0)} (w^+(\mu, v), u^+(\mu, v)) = (w^\diamond, u^\diamond).$$

Using this equality and (23), we finally conclude that

$$\begin{aligned} \lim_{(\mu, v) \rightarrow (0_+, 0)} w_{\mathcal{B}}(\mu, v) &= \lim_{(\mu, v) \rightarrow (0_+, 0)} w_{\mathcal{B}}^+(\mu, v) = w_{\mathcal{B}}^\diamond > 0, \\ \lim_{(\mu, v) \rightarrow (0_+, 0)} w_{\mathcal{N}}(\mu, v) &= \lim_{(\mu, v) \rightarrow (0_+, 0)} \mu w_{\mathcal{N}}^+(\mu, v) = 0. \end{aligned}$$

From part (b), we further see that  $(w_{\mathcal{B}}^\diamond, 0)$  is an optimal solution of (1) with maximum cardinality. This proves part (c) of the theorem.  $\blacksquare$

As an immediate consequence of Theorem 3.1, we now state a global convergence result regarding the outer iterations of our IP method.

**Corollary 3.1.** *Let  $\{\mu_k\}$  and  $\{\tilde{w}^k\}$  be the sequences generated in the primal IP method. Let  $w^k = P\tilde{w}^k + q$  for all  $k$ . Then the following statements hold:*

- (a)  $\lim_{k \rightarrow \infty} \Phi(\mathcal{M}(w^k)) = f^*$  and any accumulation point of  $\{w^k\}$  is an optimal solution of (1).
- (b) Suppose that problem (1) has an optimal solution  $w^*$  with  $\mathcal{M}(w^*) \succ 0$  and  $\epsilon(\mu_k) = O(\mu_k)$ . Then any accumulation point of  $\{w^k\}$  is an optimal solution of (1) with maximum cardinality.
- (c) Suppose that problem (1) has an optimal solution  $w^*$  with  $\mathcal{M}(w^*) \succ 0$ . Let  $u^*$  be the associated Lagrange multiplier satisfying (9). Assume that  $|\mathcal{B}| + |\mathcal{N}| = n$ , where  $\mathcal{B} := \{i : w_i^* > 0\}$  and  $\mathcal{N} := \{i : u_i^* > 0\}$ . Suppose further that  $[\nabla^2 \Phi_{\mathcal{M}}(w)]_{\mathcal{B}\mathcal{B}} \succ 0$  for any  $w \in \Omega$  satisfying  $w_{\mathcal{B}} > 0$  and  $w_{\mathcal{N}} = 0$ . Then  $\{w^k\}$  converges to an optimal solution of (1) with maximum cardinality.

*Proof.* Let  $v^k = \nabla f_{\mu_k}(\tilde{w}^k)$  and  $u^k = \mu_k(w^k)^{-1}$  for all  $k$ . In view of (10) and (11), we can observe that  $(w^k, u^k, \mu_k, v^k)$  satisfies (11). By virtue of the definition of  $(w(\mu, v), u(\mu, v))$ , we then have  $(w^k, u^k) = (w(\mu_k, v^k), u(\mu_k, v^k))$ , which together with the fact  $\mu_k \rightarrow 0$  and Theorem 3.1 implies the conclusion holds.  $\blacksquare$

Before ending this section, we establish a convergence result regarding the inner iterations of our IP method.

**Proposition 3.2.** *Let  $\mu_k > 0$  and  $\epsilon(\mu_k) > 0$  be given. Then the Newton's method detailed in step 2 of the primal IP method starting from any strictly feasible point  $\tilde{w}^{\text{init}}$  of (5) generates a point  $\tilde{w}^k$  satisfying  $\|\nabla f_{\mu_k}(\tilde{w}^k)\| \leq \epsilon(\mu_k)$  within a finite number of iterations.*

*Proof.* First, observe that all iterates generated by the Newton's method lie in the compact level set  $\Upsilon := \{\tilde{w} : f_{\mu_k}(\tilde{w}) \leq f_{\mu_k}(\tilde{w}^{\text{init}})\}$ . Furthermore, it holds that  $\tilde{w} > 0$  and  $1 - e^T \tilde{w} > 0$  for all  $\tilde{w} \in \Upsilon$ . This together with the assumption that  $\mathcal{M}(\Omega) \cap \mathcal{S}_{++}^m \neq \emptyset$  implies that  $\mathcal{M}(\Upsilon) \subset \mathcal{S}_{++}^m$ . Thus  $\nabla f_{\mu_k}$  and  $\nabla^2 f_{\mu_k}$  are continuous in  $\Upsilon$ . Using this observation and the strong convexity of  $f_{\mu_k}$  in  $\Upsilon$ , there exist  $\underline{\lambda}, \bar{\lambda} > 0$  such that  $\underline{\lambda}I \preceq \nabla^2 f_{\mu_k}(\tilde{w}) \preceq \bar{\lambda}I$  for all  $\tilde{w} \in \Upsilon$ . This relation along with the continuity

of  $\nabla f_{\mu_k}$  and  $\nabla^2 f_{\mu_k}$  implies that  $d = -(\nabla^2 f_{\mu_k}(\tilde{w}))^{-1} \nabla f_{\mu_k}(\tilde{w})$  is continuous in  $\Upsilon$ . In view of this result and the definition of  $\bar{\alpha}(\tilde{w})$ , it is not hard to show that  $\bar{\alpha}(\tilde{w})$  is positive and continuous in  $\Upsilon$ . This fact together with the compactness of  $\Upsilon$  yields  $\underline{\alpha} := \inf\{\bar{\alpha}(\tilde{w}) : \tilde{w} \in \Upsilon\} > 0$ . Thus, all iterates  $\tilde{w}$  generated by the Newton's method satisfy  $\underline{\alpha}I \preceq \nabla^2 f_{\mu_k}(\tilde{w}) \preceq \bar{\alpha}I$  and  $\bar{\alpha}(\tilde{w}) \in [\underline{\alpha}, 1]$ . The remaining proof follows the same arguments as in the proof of [22, Theorem 3.13]. ■

## 4 IP method for classical optimality criteria

In this section, we discuss how to apply our IP method to solve problem (1) with A-, D- and  $p$ th mean criterion. In particular, we will demonstrate how the Newton direction (8) can be efficiently computed for each criterion.

Before proceeding, we introduce some notations that will be used in this section (see, for example, [29] for more details). Given matrices  $A$  and  $B$  in  $\Re^{m \times n}$ ,  $A \otimes B$  denotes the Kronecker product of  $A$  and  $B$ , while  $A \circ B$  denotes the Hadamard (entry-wise) product of  $A$  and  $B$ . In addition,  $\mathbf{vec}(A)$  denotes the column vector formed by stacking columns of  $A$  one by one. For any  $m \times m$  symmetric matrix  $U$ , we define the vector  $\mathbf{svec}(U) \in \Re^{m(m+1)/2}$  as

$$\mathbf{svec}(U) = (u_{11}, \sqrt{2}u_{21}, \dots, \sqrt{2}u_{m1}, u_{22}, \sqrt{2}u_{32}, \dots, \sqrt{2}u_{m2}, \dots, u_{mm})^T.$$

It is not hard to observe that  $\mathbf{svec}$  is an isometry between  $\mathcal{S}^m$  and  $\Re^{m(m+1)/2}$  and moreover,

$$\mathrm{tr}(UV) = \mathbf{svec}(U)^T \mathbf{svec}(V) \quad \forall U, V \in \mathcal{S}^m. \quad (27)$$

We denote the inverse map of  $\mathbf{svec}$  by  $\mathbf{smat}$ . Clearly, they are adjoint of each other, namely,

$$u^T \mathbf{svec}(V) = \mathrm{tr}(\mathbf{smat}(u)V) \quad \forall u \in \Re^{m(m+1)/2}, V \in \mathcal{S}^m.$$

The symmetric Kronecker product of any two (not necessarily symmetric) matrices  $G, H \in \Re^{m \times m}$  is a square matrix of order  $m(m+1)/2$  such that

$$(G \otimes_s H) \mathbf{svec}(U) = \frac{1}{2} \mathbf{svec}(GUH^T + HUG^T) \quad \forall U \in \mathcal{S}^m. \quad (28)$$

As mentioned in [29],  $G \otimes_s H$  can be expressed in terms of the standard Kronecker product of  $G$  and  $H$  as follows:

$$G \otimes_s H = \frac{1}{2} Q(G \otimes H + H \otimes G) Q^T,$$

where  $Q \in \Re^{m(m+1)/2 \times m^2}$  is such that

$$Q \mathbf{vec}(U) = \mathbf{svec}(U), \quad Q^T \mathbf{svec}(U) = \mathbf{vec}(U) \quad \forall U \in \mathcal{S}^m. \quad (29)$$

It is easy to observe that the above  $Q$  exists and is unique. Moreover,  $QQ^T = I$ .

Throughout this section, for each optimality criterion  $\Phi$ , we define the associated function  $\phi$  as follows:

$$\phi(x) = \Phi(\mathbf{smat}(x)) \quad (30)$$

for any  $x \in \Re^{m(m+1)/2}$ , provided that  $\Phi(\mathbf{smat}(x))$  is well-defined. It is clear to observe that  $\phi$  is convex due to the convexity of  $\Phi$ . Define

$$M := [\mathbf{svec}(A_1) \dots \mathbf{svec}(A_n)].$$

Clearly,  $M \in \Re^{m(m+1)/2 \times n}$ .

With the notations above, the function  $f_\mu$  defined in (7) can be rewritten as

$$f_\mu(\tilde{w}) = \phi(M(P\tilde{w} + q)) - \mu \sum_{i=1}^{n-1} \log(\tilde{w}_i) - \mu \log(1 - e^T \tilde{w}).$$

By the chain rule, the gradient and Hessian of  $f_\mu$  are given by

$$\begin{aligned}\nabla f_\mu(\tilde{w}) &= P^T M^T \nabla \phi(Mw) - \mu P^T w^{-1}, \\ \nabla^2 f_\mu(\tilde{w}) &= P^T M^T \nabla^2 \phi(Mw) MP + \frac{\mu}{(1 - e^T \tilde{w})^2} ee^T + \mu \mathcal{D}(\tilde{w}^{-2}),\end{aligned}\quad (31)$$

where  $w = P\tilde{w} + q$ .

The main computational effort of our IP method lies in computing the Newton direction  $d$  by solving the system  $\nabla^2 f_\mu(\tilde{w})d = -\nabla f_\mu(\tilde{w})$  (see (8)). In applications,  $n$  can be significantly larger than  $m^2$ . Since the rank of  $\nabla^2 \phi(Mw)$  is at most  $m(m+1)/2$ , the first matrix in (31) has “low” rank compared to  $\nabla^2 f_\mu(\tilde{w})$ . It is generally more efficient to compute the Newton direction via the Sherman-Morrison-Woodbury formula. To this end, suppose that  $\nabla^2 \phi(Mw)$  has rank  $r$ . Let  $VDV^T$  be the *partial* eigenvalue decomposition of  $\nabla^2 \phi(Mw)$ , where  $D$  is the  $r \times r$  diagonal matrix whose diagonal consists of  $r$  largest eigenvalues of  $\nabla^2 \phi(Mw)$ , and the columns of  $V$  are the corresponding eigenvectors.<sup>1</sup> Due to the convexity of  $\phi$ , one can observe that  $\nabla^2 \phi(Mw) = VDV^T$ . It then follows from (31) that

$$\nabla^2 f_\mu(\tilde{w}) = (P^T M^T V)D(V^T MP) + \frac{\mu}{(1 - e^T \tilde{w})^2} ee^T + \mu \mathcal{D}(\tilde{w}^{-2}),$$

which together with the Sherman-Morrison-Woodbury formula yields the Newton direction

$$d = -(\nabla^2 f_\mu(\tilde{w}))^{-1} \nabla f_\mu(\tilde{w}) = -\left[ \frac{1}{\mu} \mathcal{D}(\tilde{w}^2) - \frac{1}{\mu^2} \mathcal{D}(\tilde{w}^2) (P^T M^T V \quad e) W \begin{pmatrix} V^T MP \\ e^T \end{pmatrix} \mathcal{D}(\tilde{w}^2) \right] \nabla f_\mu(\tilde{w}),$$

where

$$W = \left( \begin{pmatrix} D^{-1} & 0 \\ 0 & \frac{(1 - e^T \tilde{w})^2}{\mu} \end{pmatrix} + \frac{1}{\mu} \begin{pmatrix} V^T MP \\ e^T \end{pmatrix} \mathcal{D}(\tilde{w}^2) (P^T M^T V \quad e) \right)^{-1}.$$

When  $n \gg m^2$ , the above approach is much more efficient than solving the Newton system directly by performing Cholesky factorization of  $\nabla^2 f_\mu(\tilde{w})$ . We remark that the ideas of using Sherman-Morrison-Woodbury formula to solve specially structured Newton systems have been explored in literature (see, for example, [1, 16]).

As seen from above,  $\nabla \phi(Mw)$  and  $\nabla^2 \phi(Mw)$  are needed to compute Newton direction. For the rest of this section, we will discuss how to evaluate these two quantities for A-, D- and  $p$ th mean criterion, and determine the rank of  $\nabla^2 \phi(Mw)$  which is used to perform the aforementioned partial eigenvalue decomposition of  $\nabla^2 \phi(Mw)$ .

#### 4.1 IP method for $p$ th mean criterion

Recall from Section 1 that in  $\mathcal{S}_{++}^m$ , the  $p$ th mean criterion  $\Phi$  becomes

$$\Phi(X) = \text{tr}((K^T X^{-1} K)^{-p}) \quad (32)$$

for some  $p < 0$  and  $K \in \mathbb{R}^{m \times k}$  with full column rank. It is easy to check that Assumption 3.1 holds for  $\Phi$ . Hence, problem (1) with this criterion can be suitably solved by our IP method proposed in Section 3.

Based on the above discussion, we know that our IP method needs the gradient and Hessian of the associated function  $\phi$  for computing Newton direction, where  $\phi$  is defined by (30). We next discuss how to compute them. Before proceeding, we state the following classical result (see, for example, [11, Proposition 4.3]) that will be used subsequently.

---

<sup>1</sup>The partial eigenvalue decomposition can be efficiently computed by the package PROPACK [21].

**Lemma 4.1.** Let  $g : \mathbb{R} \rightarrow \mathbb{R}$  be a differentiable function and let  $g^\square : \mathcal{S}^m \rightarrow \mathcal{S}^m$  be defined by

$$g^\square(Y) := V \begin{pmatrix} g(d_1) & & & \\ & g(d_2) & & \\ & & \ddots & \\ & & & g(d_m) \end{pmatrix} V^T,$$

where  $V\mathcal{D}(d)V^T$  is an eigenvalue decomposition of  $Y$  for some  $d \in \mathbb{R}^m$ . Then the function  $g^\square$  is well-defined, i.e., it is independent of the choice of  $V$  and  $d$ , and is also differentiable. Moreover, let  $S^{g,d} \in \mathcal{S}^m$  be a symmetric matrix whose  $(i, j)$ th entry is given by

$$s_{ij}^{g,d} := \begin{cases} \frac{g(d_i) - g(d_j)}{d_i - d_j} & \text{if } d_i \neq d_j, \\ g'(d_i) & \text{otherwise.} \end{cases}$$

Then the directional derivative of  $g^\square$  at  $Y$  along the direction  $H \in \mathcal{S}^m$  is given by

$$V(S^{g,d} \circ (V^T H V))V^T.$$

**Proposition 4.1.** Let  $\Phi$  be defined in (32) and the associated  $\phi$  be defined in (30). Let  $Q \in \mathbb{R}^{m(m+1)/2 \times m^2}$  be defined in (29). Then the gradient and Hessian of  $\phi$  at any  $x \in \mathbf{svec}(\mathcal{S}_{++}^m)$  are given by

$$\nabla \phi(x) = p \mathbf{svec}(X^{-1}K(K^T X^{-1}K)^{-p-1}K^T X^{-1}), \quad (33)$$

$$\begin{aligned} \nabla^2 \phi(x) = & Q(-p[(X^{-1}KV) \otimes (X^{-1}KV)]\mathcal{D}(\mathbf{vec}(S^{g,d}))[(X^{-1}KV) \otimes (X^{-1}KV)]^T \\ & - p X^{-1} \otimes G - p G \otimes X^{-1})Q^T, \end{aligned} \quad (34)$$

respectively, where  $X = \mathbf{smat}(x)$ ,  $V\mathcal{D}(d)V^T$  is an eigenvalue decomposition of  $K^T X^{-1}K$  for some  $d \in \mathbb{R}^m$ ,  $g(t) = t^{-p-1}$ , and  $G = X^{-1}K[K^T X^{-1}K]^{-p-1}K^T X^{-1}$ . In particular, when  $K = I$ , the above gradient and Hessian reduce to

$$\nabla \phi(x) = p \mathbf{svec}(X^{p-1}), \quad (35)$$

$$\nabla^2 \phi(x) = Q(V \otimes V)\mathcal{D}(\mathbf{vec}(S^{g,d}))(V \otimes V)^T Q^T, \quad (36)$$

where  $g(t) = pt^{p-1}$  and  $V\mathcal{D}(d)V^T$  is an eigenvalue decomposition of  $X$  for some  $d \in \mathbb{R}^m$ .

*Proof.* To derive the gradient of  $\phi$ , we fix an arbitrary  $x \in \mathbf{svec}(\mathcal{S}_{++}^m)$ . Let  $X = \mathbf{smat}(x)$ . For all sufficiently small  $h \in \mathbb{R}^{m(m+1)/2}$ , we have  $X + H \succ 0$ , where  $H = \mathbf{smat}(h)$ , and moreover,

$$(X + H)^{-1} = X^{-1} - X^{-1}HX^{-1} + o(H). \quad (37)$$

Using (37) and Lemma 4.1 with  $g(t) = t^{-p}$  and  $Y = K^T X^{-1}K$ , we obtain that

$$\begin{aligned} \Phi(X + H) &= \text{tr}((K^T[X + H]^{-1}K)^{-p}) = \text{tr}((K^T X^{-1}K - K^T X^{-1}HX^{-1}K + o(H))^{-p}) \\ &= \Phi(X) - \text{tr}(V(S^{g,d} \circ (V^T K^T X^{-1}HX^{-1}KV))V^T) + o(H), \end{aligned} \quad (38)$$

where  $V\mathcal{D}(d)V^T$  is an eigenvalue decomposition of  $Y$ . Letting  $R := -K^T X^{-1}HX^{-1}K$  and using the fact that  $V^T V = I$  and  $s_{ii}^{g,d} = -pd_i^{-p-1}$  for all  $i$ , we further have

$$\begin{aligned} \text{tr}(V(S^{g,d} \circ (V^T R V))V^T) &= \text{tr}(S^{g,d} \circ (V^T R V)) = \sum_{i=1}^m s_{ii}^{g,d} \sum_{j,k} v_{ji} r_{jk} v_{ki} \\ &= -p \sum_{j,k} \left( \sum_{i=1}^m v_{ji} d_i^{-p-1} v_{ki} \right) r_{jk} = -\text{tr}(p(K^T X^{-1}K)^{-p-1}R) \\ &= \text{tr}(pX^{-1}K(K^T X^{-1}K)^{-p-1}K^T X^{-1}H). \end{aligned} \quad (39)$$

In view of the definitions of  $\phi$ ,  $\Phi$ ,  $X$  and  $H$ , it follows from (38), (39) and (27) that

$$\phi(x+h) - \phi(x) = \Phi(X+H) - \Phi(X) = h^T (p \mathbf{svec}(X^{-1}K(K^T X^{-1}K)^{-p-1}K^T X^{-1})) + o(h),$$

which yields (33). And (35) immediately follows from (33) by letting  $K = I$ .

We next derive the Hessian of  $\phi$  at any  $x \in \mathbf{svec}(\mathcal{S}_{++}^m)$ . To proceed, we first recall the following well-known results (see, for example, page 243 and Lemma 4.3.1 of [20]):

$$\mathbf{vec}(ABC) = (C^T \otimes A) \mathbf{vec}(B), \quad (A \otimes B)^T = A^T \otimes B^T. \quad (40)$$

Let  $X$ ,  $h$  and  $H$  be defined as above. Using (37) and Lemma 4.1 with  $g(t) = t^{-p-1}$  and  $Y = K^T X^{-1}K$ , we have

$$\begin{aligned} \nabla \Phi(X+H) &= p(X+H)^{-1}K[K^T(X+H)^{-1}K]^{-p-1}K^T(X+H)^{-1} \\ &= p(X^{-1} - X^{-1}HX^{-1})K[K^T(X^{-1} - X^{-1}HX^{-1})K]^{-p-1}K^T(X^{-1} - X^{-1}HX^{-1}) + o(H) \\ &= \nabla \Phi(X) - p(X^{-1}K)V(S^{g,d} \circ (V^T K^T X^{-1}HX^{-1}KV))V^T(K^T X^{-1}) \\ &\quad - pGHX^{-1} - pX^{-1}HG + o(H), \end{aligned} \quad (41)$$

where  $G$  is defined as above. Since  $X$  is symmetric, it follows from (40) that

$$\begin{aligned} &\mathbf{vec}((X^{-1}KV)(S^{g,d} \circ (V^T K^T X^{-1}HX^{-1}KV))(V^T K^T X^{-1})) \\ &= [(X^{-1}KV) \otimes (X^{-1}KV)] \mathbf{vec}(S^{g,d} \circ (V^T K^T X^{-1}HX^{-1}KV)) \\ &= [(X^{-1}KV) \otimes (X^{-1}KV)] \mathcal{D}(\mathbf{vec}(S^{g,d})) \mathbf{vec}(V^T K^T X^{-1}HX^{-1}KV) \\ &= [(X^{-1}KV) \otimes (X^{-1}KV)] \mathcal{D}(\mathbf{vec}(S^{g,d})) [(X^{-1}KV) \otimes (X^{-1}KV)]^T \mathbf{vec}(H). \end{aligned} \quad (42)$$

In addition, since  $G$  is symmetric, we further have that

$$\mathbf{vec}(GHX^{-1} + X^{-1}HG) = [X^{-1} \otimes G + G \otimes X^{-1}] \mathbf{vec}(H). \quad (43)$$

In addition, by virtue of (29), (30), the definition of  $X$  and  $H$ , and the fact that  $\mathbf{svec}$  is the adjoint operator of  $\mathbf{smat}$ , one can have

$$\nabla \phi(x+h) - \nabla \phi(x) = \mathbf{svec}(\nabla \Phi(X+H) - \nabla \Phi(X)) = Q \mathbf{vec}(\nabla \Phi(X+H) - \nabla \Phi(X)).$$

This relation together with (29), (41)–(43) and the definition of  $H$  yields

$$\begin{aligned} \nabla \phi(x+h) - \nabla \phi(x) &= Q(-p[(X^{-1}KV) \otimes (X^{-1}KV)] \mathcal{D}(\mathbf{vec}(S^{g,d}))[(X^{-1}KV) \otimes (X^{-1}KV)]^T \\ &\quad - pX^{-1} \otimes G - pG \otimes X^{-1}) \mathbf{vec}(H) + o(Q \mathbf{vec}(H)) \\ &= Q(-p[(X^{-1}KV) \otimes (X^{-1}KV)] \mathcal{D}(\mathbf{vec}(S^{g,d}))[(X^{-1}KV) \otimes (X^{-1}KV)]^T \\ &\quad - pX^{-1} \otimes G - pG \otimes X^{-1}) Q^T \mathbf{svec}(H) + o(\mathbf{svec}(H)) \\ &= Q(-p[(X^{-1}KV) \otimes (X^{-1}KV)] \mathcal{D}(\mathbf{vec}(S^{g,d}))[(X^{-1}KV) \otimes (X^{-1}KV)]^T \\ &\quad - pX^{-1} \otimes G - pG \otimes X^{-1}) Q^T h + o(h), \end{aligned}$$

and hence (34) holds.

For the case when  $K = I$ ,  $\nabla^2 \phi$  can be directly derived as follows. We know from (35) that  $\nabla \Phi(X) = pX^{p-1}$ . Letting  $g(t) = pt^{p-1}$  and  $V\mathcal{D}(d)V^T$  be an eigenvalue decomposition of  $X$ , it follows from Lemma 4.1 that

$$\nabla \Phi(X+H) = \nabla \Phi(X) + V(S^{g,d} \circ (V^T HV))V^T + o(H).$$

In view of (40), one can see that

$$\begin{aligned} \mathbf{vec}(V(S^{g,d} \circ (V^T HV))V^T) &= (V \otimes V) \mathbf{vec}(S^{g,d} \circ (V^T HV)) \\ &= (V \otimes V)[\mathbf{vec}(S^{g,d}) \circ \mathbf{vec}(V^T HV)] \\ &= (V \otimes V)(\mathbf{vec}(S^{g,d}) \circ [(V \otimes V)^T \mathbf{vec}(H)]) \\ &= (V \otimes V) \mathcal{D}(\mathbf{vec}(S^{g,d}))(V \otimes V)^T \mathbf{vec}(H). \end{aligned}$$

Using these relations and a similar proof as above, we can see that (36) holds.  $\blacksquare$

As mentioned earlier, we need to know the rank of  $\nabla^2\phi(x)$  for performing the partial eigenvalue decomposition of  $\nabla^2\phi(x)$  which is used to compute Newton direction. In the next proposition, we determine the rank of  $\nabla^2\phi(x)$  at any  $x \in \mathbf{svec}(\mathcal{S}_{++}^m)$ .

**Proposition 4.2.** *Let  $\Phi$  be defined in (32) and the associated  $\phi$  be defined in (30). Then the rank of  $\nabla^2\phi(x)$  is  $m(m+1)/2 - (m-k)(m-k+1)/2$  for any  $x \in \mathbf{svec}(\mathcal{S}_{++}^m)$ .*

*Proof.* Let  $x \in \mathbf{svec}(\mathcal{S}_{++}^m)$  be arbitrarily chosen. Define  $X = \mathbf{smat}(x)$ . Let  $G$ ,  $V$ ,  $d$  and  $S^{g,d}$  be defined in Proposition 4.1 with  $g(t) = t^{-p-1}$ . For convenience, we define

$$\begin{aligned} M_1 &= [(X^{-1}KV) \otimes (X^{-1}KV)] \mathcal{D}(\mathbf{vec}(S^{g,d})) [(X^{-1}KV) \otimes (X^{-1}KV)]^T, \\ M_2 &= X^{-1} \otimes G + G \otimes X^{-1}. \end{aligned}$$

To determine the rank of  $\nabla^2\phi(x)$ , it suffices to know the dimension of the null space of  $\nabla^2\phi(x)$ , denoted by  $\text{Null}(\nabla^2\phi(x))$ . Notice that  $\phi$  is a twice differentiable convex function in  $\mathbf{svec}(\mathcal{S}_{++}^m)$ . Thus,  $\nabla^2\phi(x) \succeq 0$ . It implies that  $h \in \text{Null}(\nabla^2\phi(x))$  if and only if  $h^T \nabla^2\phi(x) h = 0$ . We will subsequently show that

$$h^T \nabla^2\phi(x) h = 0 \Leftrightarrow K^T X^{-1} H = 0, \quad (44)$$

where  $H = \mathbf{smat}(h)$ . It then follows that

$$h \in \text{Null}(\nabla^2\phi(x)) \Leftrightarrow K^T X^{-1} H = 0.$$

Notice that  $K^T X^{-1}$  has full row rank. Thus, there exist nonsingular matrices  $E_1$  and  $E_2$  such that  $K^T X^{-1} = E_1 \begin{pmatrix} I & 0 \end{pmatrix} E_2$ , where  $I$  is the identity matrix of order  $k$ . It then follows that

$$K^T X^{-1} H = 0 \Leftrightarrow \begin{pmatrix} I & 0 \end{pmatrix} U = 0,$$

where  $U = E_2 H E_1^T \in \mathcal{S}^m$ . It is easy to see that the dimension of  $\{U \in \mathcal{S}^m : \begin{pmatrix} I & 0 \end{pmatrix} U = 0\}$  is  $(m-k)(m-k+1)/2$ . Since  $E_2$  is invertible, we conclude that the dimension of  $\{H \in \mathcal{S}^m : K^T X^{-1} H = 0\}$  is also  $(m-k)(m-k+1)/2$ . Since  $\mathbf{smat}$  is a one-to-one map between  $\mathfrak{R}^{m(m+1)/2}$  and  $\mathcal{S}^m$ , the dimension of  $\text{Null}(\nabla^2\phi(x))$  is  $(m-k)(m-k+1)/2$ , and hence the rank of  $\nabla^2\phi(x)$  is  $m(m+1)/2 - (m-k)(m-k+1)/2$ . To complete the proof, we next show that (44) holds by considering two cases  $p \leq -1$  or  $-1 < p < 0$ .

We start with the first case  $p \leq -1$ . Notice that all entries of  $S^{g,d}$  are nonnegative and thus  $M_1 \succeq 0$ . Also,  $M_2 \succeq 0$ . It then follows from (34) and (29) that  $h^T \nabla^2\phi(x) h = 0$  if and only if

$$\mathbf{vec}(H)^T M_1 \mathbf{vec}(H) = 0, \quad \mathbf{vec}(H)^T M_2 \mathbf{vec}(H) = 0. \quad (45)$$

By (43), the second equality of (45) becomes

$$\text{tr}(H X^{-1} K (K^T X^{-1} K)^{-p-1} K^T X^{-1} H X^{-1}) = 0,$$

which is equivalent to

$$\begin{aligned} \text{tr}(X^{-\frac{1}{2}} H X^{-1} K (K^T X^{-1} K)^{-\frac{p+1}{2}} (K^T X^{-1} K)^{-\frac{p+1}{2}} K^T X^{-1} H X^{-\frac{1}{2}}) &= 0 \\ \Leftrightarrow (K^T X^{-1} K)^{-\frac{p+1}{2}} K^T X^{-1} H X^{-\frac{1}{2}} &= 0 \Leftrightarrow K^T X^{-1} H = 0. \end{aligned} \quad (46)$$

Moreover,  $K^T X^{-1} H = 0$  implies that the first equality of (45) holds. Therefore, (45) holds if and only if  $K^T X^{-1} H = 0$ . It follows that (44) holds for  $p \leq -1$ .

We next show that (44) also holds for  $-1 < p < 0$ . Indeed, for such  $p$ , all entries of  $S^{g,d}$  are negative and hence  $-M_1 \succeq 0$ . Using Proposition 4.1, we see that  $h^T \nabla^2\phi(x) h = 0$  if and only if

$$\mathbf{vec}(H)^T (M_1 + M_2) \mathbf{vec}(H) = 0. \quad (47)$$



We claim that

$$\frac{1}{2} \mathbf{vec}(H)^T M_2 \mathbf{vec}(H) \geq -\mathbf{vec}(H)^T M_1 \mathbf{vec}(H). \quad (48)$$

Indeed, letting  $W = (K^T X^{-1} K)^{-1}$  and using Lemma 3.1, we have

$$W^{-1} = K^T X^{-1} K \Rightarrow \begin{pmatrix} X & K \\ K^T & W^{-1} \end{pmatrix} \succeq 0 \Rightarrow X \succeq K W K^T.$$

The latter relation together with the definitions of  $M_2$ ,  $G$  and (43) implies that

$$\begin{aligned} \frac{1}{2} \mathbf{vec}(H)^T M_2 \mathbf{vec}(H) &= \text{tr}(H X^{-1} H X^{-1} K W^{p+1} K^T X^{-1}) \\ &= \text{tr}([X^{-1} H (X^{-1} K W^{p+1} K^T X^{-1})^{\frac{1}{2}}]^T X [X^{-1} H (X^{-1} K W^{p+1} K^T X^{-1})^{\frac{1}{2}}]) \\ &\geq \text{tr}([X^{-1} H (X^{-1} K W^{p+1} K^T X^{-1})^{\frac{1}{2}}]^T K W K^T [X^{-1} H (X^{-1} K W^{p+1} K^T X^{-1})^{\frac{1}{2}}]) \\ &= \text{tr}(H X^{-1} K W K^T X^{-1} H X^{-1} K W^{p+1} K^T X^{-1}) \end{aligned} \quad (49)$$

Let  $Z = V^T K^T X^{-1} H X^{-1} K V$ . Notice that  $W = V \mathcal{D}(d^{-1}) V^T$ . Using this relation, the definition of  $Z$  and (40), we have

$$\begin{aligned} \text{tr}(H X^{-1} K W K^T X^{-1} H X^{-1} K W^{p+1} K^T X^{-1}) &= \text{tr}(H X^{-1} K V \mathcal{D}(d^{-1}) Z \mathcal{D}(d^{-p-1}) V^T K^T X^{-1}) \\ &= \text{tr}(\mathcal{D}(d^{-1}) Z \mathcal{D}(d^{-p-1}) Z) = \mathbf{vec}(Z)^T [\mathcal{D}(d^{-1}) \otimes \mathcal{D}(d^{-p-1})] \mathbf{vec}(Z), \end{aligned}$$

which together with (49) yields

$$\frac{1}{2} \mathbf{vec}(H)^T M_2 \mathbf{vec}(H) \geq \mathbf{vec}(Z)^T [\mathcal{D}(d^{-1}) \otimes \mathcal{D}(d^{-p-1})] \mathbf{vec}(Z). \quad (50)$$

Also, by the definitions of  $M_1$ ,  $Z$  and (40), we obtain that

$$-\mathbf{vec}(H)^T M_1 \mathbf{vec}(H) = \mathbf{vec}(Z)^T \mathcal{D}(\mathbf{vec}(-S^{g,d})) \mathbf{vec}(Z). \quad (51)$$

Since  $1 > p + 1 > 0$  and  $d_i > 0$  for all  $i$ , it is not hard to show that

$$d_i^{-1} d_j^{-p-1} \geq -\frac{d_i^{-p-1} - d_j^{-p-1}}{d_i - d_j},$$

whenever  $d_i \neq d_j$ . Thus,  $\mathcal{D}(d^{-1}) \otimes \mathcal{D}(d^{-p-1}) \succeq \mathcal{D}(\mathbf{vec}(-S^{g,d}))$ , which together with (50) and (51) implies that (48) holds. It then follows from (48), (47) and the fact  $M_2 \succeq 0$  that  $\mathbf{vec}(H)^T M_2 \mathbf{vec}(H) = 0$ . The rest of proof is similar to the case  $p \leq -1$ . ■

## 4.2 IP method for A-criterion

Recall from Section 1 that in  $\mathcal{S}_{++}^m$ , the A-criterion  $\Phi$  becomes

$$\Phi(X) = \text{tr}(K^T X^{-1} K) \quad (52)$$

for some  $K \in \mathbb{R}^{m \times k}$  with full column rank. Since A-criterion is a special case of  $p$ th mean criterion, the IP method discussed in Sections 3 and 4.1 can be suitably applied to solve problem (1) with A-criterion. We next show that by exploiting the special structure, we can obtain a more compact representation of the associated Hessian matrix that is used to compute Newton direction for our IP method.

**Proposition 4.3.** *Let  $\Phi$  be defined in (52) and the associated  $\phi$  be defined in (30). Then the gradient and Hessian of  $\phi$  at any  $x \in \mathbf{svec}(\mathcal{S}_{++}^m)$  are given by*

$$\nabla \phi(x) = -\mathbf{svec}(X^{-1} K K^T X^{-1}), \quad (53)$$

$$\nabla^2 \phi(x) = 2X^{-1} \otimes_s (X^{-1} K K^T X^{-1}), \quad (54)$$

where  $X = \mathbf{smat}(x)$ .

*Proof.* (53) follows immediately from (33) with  $p = -1$ . We now prove (54). Let  $x \in \mathbf{svec}(\mathcal{S}_{++}^m)$  be arbitrarily chosen, and let  $X = \mathbf{smat}(x)$ . For all sufficiently small  $h \in \Re^{m(m+1)/2}$ , we observe  $X + H \succ 0$ , where  $H = \mathbf{smat}(h)$ . In view of the definitions of  $X$  and  $H$ , it then follows from (53), (37) and (28) that

$$\begin{aligned}\nabla\phi(x+h) - \nabla\phi(x) &= -\mathbf{svec}\left((X+H)^{-1}KK^T(X+H)^{-1} - X^{-1}KK^TX^{-1}\right) \\ &= \mathbf{svec}(X^{-1}HX^{-1}KK^TX^{-1} + X^{-1}KK^TX^{-1}HX^{-1}) + o(\mathbf{svec}(H)) \\ &= 2X^{-1}\otimes_s(X^{-1}KK^TX^{-1})h + o(h),\end{aligned}$$

which proves (54).  $\blacksquare$

Since the A-criterion is a special case of the  $p$ th mean criterion, it follows from Proposition 4.2 that the rank of  $\nabla^2\phi(x)$  is also  $m(m+1)/2 - (m-k)(m-k+1)/2$  for every  $x \in \mathbf{svec}(\mathcal{S}_{++}^m)$ .

### 4.3 IP method for D-criterion

Recall from Section 1 that in  $\mathcal{S}_{++}^m$ , the D-criterion  $\Phi$  becomes

$$\Phi(X) = \log \det(K^TX^{-1}K) \quad (55)$$

for some  $K \in \Re^{m \times k}$  with full column rank. It is easy to verify that Assumption 3.1 is satisfied. Hence, problem (1) with this criterion can be suitably solved by the IP method studied in Section 3. In the next proposition, we provide formulas for computing gradient and Hessian of the associated function  $\phi$  that are used in the IP method. The proof is similar to that of Proposition 4.3 and is thus omitted.

**Proposition 4.4.** *Let  $\Phi$  be defined in (55) and the associated  $\phi$  be defined in (30). Then the gradient and Hessian of  $\phi$  at any  $x \in \mathbf{svec}(\mathcal{S}_{++}^m)$  are given by*

$$\begin{aligned}\nabla\phi(x) &= -\mathbf{svec}(X^{-1}KWK^TX^{-1}), \\ \nabla^2\phi(x) &= 2X^{-1}\otimes_s(X^{-1}KWK^TX^{-1}) - (X^{-1}KWK^TX^{-1})\otimes_s(X^{-1}KWK^TX^{-1}),\end{aligned} \quad (56)$$

where  $X = \mathbf{smat}(x)$  and  $W = (K^TX^{-1}K)^{-1}$ .

We next determine the rank of  $\nabla^2\phi(X)$  at any  $x \in \mathbf{svec}(\mathcal{S}_{++}^m)$ .

**Proposition 4.5.** *Let  $\Phi$  be defined in (55) and the associated  $\phi$  be defined in (30). Then the rank of  $\nabla^2\phi(x)$  is  $m(m+1)/2 - (m-k)(m-k+1)/2$  for any  $x \in \mathbf{svec}(\mathcal{S}_{++}^m)$ .*

*Proof.* Let  $x \in \mathbf{svec}(\mathcal{S}_{++}^m)$  be arbitrarily chosen. Define  $X = \mathbf{smat}(x)$ . As in the proof of Proposition 4.2, to determine the rank of  $\nabla^2\phi(x)$ , it suffices to know the dimension of  $\text{Null}(\nabla^2\phi(x))$ . Notice that  $\phi$  is a twice differentiable convex function in  $\mathbf{svec}(\mathcal{S}_{++}^m)$ . Thus,  $\nabla^2\phi(x) \succeq 0$ . It implies that  $h \in \text{Null}(\nabla^2\phi(x))$  if and only if  $h^T\nabla^2\phi(x)h = 0$ . In view of (28) and (56), it is not hard to verify that  $h^T\nabla^2\phi(x)h = 0$  if and only if

$$2\text{tr}(HX^{-1}HX^{-1}KWK^TX^{-1}) - \text{tr}(HX^{-1}KWK^TX^{-1}HX^{-1}KWK^TX^{-1}) = 0,$$

where  $H = \mathbf{smat}(h)$ . In addition, we can observe that (49) also holds for  $p = 0$ , and hence

$$\text{tr}(HX^{-1}HX^{-1}KWK^TX^{-1}) \geq \text{tr}(HX^{-1}KWK^TX^{-1}HX^{-1}KWK^TX^{-1}).$$

Furthermore,

$$\text{tr}(HX^{-1}HX^{-1}KWK^TX^{-1}) = \text{tr}\left([(X^{-1}KWK^TX^{-1})^{\frac{1}{2}}H]X^{-1}[H(X^{-1}KWK^TX^{-1})^{\frac{1}{2}}]\right) \geq 0.$$

The above relations imply that  $h^T\nabla^2\phi(x)h = 0$  if and only if

$$\text{tr}(HX^{-1}HX^{-1}KWK^TX^{-1}) = 0,$$

which together with definition of  $W$  and the same arguments used in (46) implies that

$$h^T\nabla^2\phi(x)h = 0 \Leftrightarrow K^TX^{-1}H = 0.$$

The rest of the proof follows similarly as that of Proposition 4.2.  $\blacksquare$

## 5 Computational results

In this section, we conduct numerical experiments to test the performance of the IP method discussed in this paper for solving problem (1) with A-, D- and  $p$ th mean criterion and also compare its performance with the multiplicative algorithm.

We develop Matlab codes for our IP method to solve (1) with A-, D- and  $p$ th mean criterion. We also implement the multiplicative algorithm in Matlab for solving (1) with A-, D- and  $p$ th mean criterion. To benchmark the performance of our IP method, we also report the computational results using a general SDP solver, namely, SDPT3 [30, 34] (Version 4.0) on solving a linear SDP reformulation of (1) with A-criterion (see [14, Page 532]) and a log-determinant SDP reformulation of (1) with D-criterion (see [23, Equation (9)]). We shall mention that it is not clear whether problem (1) with  $p$ th mean criterion can be reformulated into a problem that can be efficiently solved by SDPT3. As SDPT3 implements an infeasible path-following algorithm, we project the approximate solution  $w$  found by SDPT3 onto the unit simplex to obtain an approximate optimal feasible solution for problem (1) and the final objective value reported in our tests is based on the latter solution. All computations in this section are performed in Matlab 7.13.0 (2011b) on a workstation with an Intel Xeon E5410 CPU (2.33 GHz) and 8GB RAM running Red Hat Enterprise Linux (kernel 2.6.18).

For our IP method, we set  $\tilde{w}^0 = \frac{1}{n}e \in \mathbb{R}^{n-1}$ ,  $\mu_1 = 10$ ,  $\beta = \gamma = 0.5$ ,  $\sigma = 0.1$  and  $\eta = 0.95$ . In addition, we set  $\epsilon(\mu) = \max\{\mu, 1e - 10\}$  and terminate the algorithm once  $\mu_k \leq 1e - 10$ . On the other hand, for the multiplicative algorithm, similarly as in [37], we set  $\lambda = 1$ ,  $w^0 = \frac{1}{n}e \in \mathbb{R}^n$ , and terminate the algorithm when it reaches 10000 iterations or

$$\max_{1 \leq i \leq n} d_i(w^k) \leq (1 + \delta) \sum_{i=1}^n w_i^k d_i(w^k)$$

holds with  $\delta = 2e - 4$ , where  $d_i(w)$  is defined in (2). Furthermore, for SDPT3, we use the default tolerance. Finally, we use the mex files skron, smat and svec from the SDPT3 package for efficient operations on symmetric matrices in our implementation of the IP method and the multiplicative algorithm.

In our tests below, we consider the following four design spaces:

$$\begin{aligned} \chi_1(n) &= \{x_i = (e^{-s_i}, s_i e^{-s_i}, e^{-2s_i}, s_i e^{-2s_i})^T, 1 \leq i \leq n\}, \\ \chi_2(n) &= \{x_i = (1, s_i, s_i^2, s_i^3)^T, 1 \leq i \leq n\}, \\ \chi_3(n) &= \{x_{(i-1)n+j} = (1, r_i, r_i^2, t_j, r_i t_j)^T, 1 \leq i, j \leq n\}, \\ \chi_4(n) &= \{x_i = (t_i, t_i^2, \sin(2\pi t_i), \cos(2\pi t_i))^T, 1 \leq i \leq n\}, \end{aligned}$$

where  $s_i = \frac{3i}{n}$ ,  $r_i = \frac{2i}{n} - 1$  and  $t_i = \frac{i}{n}$ . The space  $\chi_1(n)$  represents the linearization of a compartmental model [3]. The space  $\chi_2(n)$  corresponds to polynomial regression. The third space, as described in [38], represents a response surface with a nonlinear effect and an interaction, while the fourth space is the quadratic/trigonometric example proposed in [36]. The test sets  $\chi_1$ ,  $\chi_3$  and a variant of the test set  $\chi_2$  are also used in [38].

In our first test, for each design space, we set  $A_i = x_i x_i^T$  for  $i = 1, \dots, n$ , with  $n = 10000, 50000, 100000$ . For each  $n$  and each design space, we randomly generate 30 different matrices  $K \in \mathbb{R}^{m \times 3}$  (i.e., we set  $k = 3$ ), each having i.i.d. Gaussian entries of mean 0 and variance 1. We then apply our IP method and the multiplicative algorithm to solve problem (1) with A-, D- and  $p$ th mean criterion on these instances and also apply SDPT3 to solve (1) with A- and D-criterion. The computational results averaged over the 30 instances are reported in Tables 1–4. In particular, the performance of our IP method, the multiplicative algorithm and SDPT3 are reported under the columns named “IP”, “MUL” and “SDPT3”, respectively. In addition, the CPU time abbreviated as “cpu” is in seconds and the objective value abbreviated as “obj” is rounded off to six significant digits. We see that our IP method significantly outperforms the multiplicative algorithm in terms of CPU time, and gives a smaller objective value in all instances. Moreover, our IP method also outperforms SDPT3 in CPU time and gives a smaller objective value

Table 1: Computational results for A-criterion with random  $K$ 

$\chi_i$	$n$	cpu			obj		
		MUL	IP	SDPT3	MUL	IP	SDPT3
1	10000	14.30	0.80	2.19	111511	110530	113966
1	50000	67.75	4.05	10.92	120146	119140	133893
1	100000	142.77	8.17	18.76	187986	186276	215851
2	10000	17.41	0.84	1.98	206.959	203.762	203.762
2	50000	78.66	4.20	10.45	269.029	264.054	264.055
2	100000	177.25	8.19	22.26	256.361	251.375	251.391
3	100	41.15	1.08	2.30	42.1423	42.1203	42.1203
3	200	139.58	4.97	14.24	54.7132	54.6797	54.6879
3	300	351.40	11.71	33.84	54.0963	54.078	54.0817
4	10000	13.02	0.94	1.81	560.199	545.569	545.569
4	50000	61.18	4.77	9.84	466.661	455.301	455.301
4	100000	138.10	9.19	20.73	611.606	597.149	597.175

Table 2: Computational results for D-criterion with random  $K$ 

$\chi_i$	$n$	cpu			obj		
		MUL	IP	SDPT3	MUL	IP	SDPT3
1	10000	1.26	0.97	1.46	19.7352	19.7347	19.7356
1	50000	5.13	4.74	5.96	19.9312	19.9307	19.933
1	100000	14.20	8.72	12.29	19.7973	19.7968	19.7987
2	10000	1.95	0.77	1.49	5.95269	5.95229	5.9523
2	50000	20.88	3.91	6.29	5.30436	5.3039	5.3039
2	100000	53.35	7.85	13.26	5.08652	5.08608	5.08609
3	100	4.18	1.05	1.76	6.58713	6.58694	6.58694
3	200	19.01	4.37	8.80	6.65124	6.65104	6.65103
3	300	68.40	9.91	21.89	6.74346	6.74327	6.74382
4	10000	1.45	0.96	1.33	7.40587	7.40535	7.40535
4	50000	13.47	4.51	5.96	7.65401	7.6535	7.6535
4	100000	40.19	8.33	12.08	8.66619	8.66575	8.66574

in most instances. Furthermore, it is worth pointing out that SDPT3 early terminates when solving some instances for  $\chi_1$  with A-criterion, possibly due to bad scaling of  $\mathcal{M}(w)$ . This accounts for its significantly larger objective values in Table 1 corresponding to  $\chi_1$ . Finally, for  $p$ th mean criterion with  $p < -1$ , our IP method achieves significantly better objective values than the multiplicative algorithm, where the objective value of the latter algorithm is chosen to be the minimum over all iterations (see Table 4). This phenomenon is actually not surprising since the multiplicative algorithm is only known to converge for  $p \in (-1, 0)$ , but it may not converge when  $p < -1$ .

In our second test, we consider the case when  $K = I$ . The instances used in this test are the same as those in the first test except  $K = I$ . We also apply our IP method and the multiplicative algorithm to solve problem (1) with A-, D- and  $p$ th mean criterion on these instances and apply SDPT3 to solve (1) with A- and D-criterion. The computational results are reported in Tables 5–8. We again observe that our IP method outperforms the multiplicative algorithm in terms of objective value in all instances, and is generally much faster on large instances. Furthermore, our IP method is usually faster than SDPT3 and produces comparable or smaller objective values.

## 6 Concluding remarks

In this paper we propose a primal IP method for solving problem (1) with a large class of convex optimality criteria and establish its global convergence. We demonstrate how the Newton direction

Table 3: Computational results for  $p$ th mean criterion with random  $K$  for some  $p \in (-1, 0)$ 

$\chi_i$	$n$	$p = -0.25$				$p = -0.75$			
		cpu		obj		cpu		obj	
		MUL	IP	MUL	IP	MUL	IP	MUL	IP
1	10000	6.43	0.90	25.4567	25.4558	9.22	0.84	5187.73	5187.23
1	50000	38.46	4.00	25.1902	25.1894	48.09	3.83	7128.51	7126.32
1	100000	87.55	7.84	25.1312	25.1304	126.19	7.69	7207.59	7205.68
2	10000	6.24	0.81	5.68067	5.68046	13.76	0.86	46.4144	46.4008
2	50000	28.79	3.71	6.00911	6.00886	75.97	4.02	58.2028	58.1903
2	100000	74.43	7.56	6.12458	6.12434	159.49	7.80	60.9256	60.9108
3	100	5.45	0.97	5.58387	5.58379	3.65	1.04	24.8691	24.868
3	200	19.26	4.23	5.56727	5.56718	17.44	4.88	23.5463	23.5451
3	300	60.36	10.64	5.45907	5.45899	58.36	11.92	24.7679	24.7664
4	10000	1.71	0.94	7.30484	7.30456	2.73	0.96	118.871	118.859
4	50000	8.81	4.57	7.27622	7.27589	10.92	4.86	108.079	108.066
4	100000	32.17	9.54	7.30129	7.30102	46.79	10.30	128.676	128.662

Table 4: Computational results for  $p$ th mean criterion with random  $K$  for some  $p < -1$ 

$\chi_i$	$n$	$p = -1.1$				$p = -1.2$			
		cpu		obj		cpu		obj	
		mul	IP	mul	IP	mul	IP	mul	IP
1	10000	4.54	0.73	611960	602294	4.12	0.71	1.46813e+06	1.43891e+06
1	50000	19.79	3.45	541355	532904	18.54	3.43	2.0649e+06	2.02777e+06
1	100000	49.28	7.07	371942	365201	52.07	6.89	1.79042e+06	1.75803e+06
2	10000	5.53	0.85	373.376	359.802	4.69	0.84	650.345	629.288
2	50000	20.36	3.95	492.463	476.123	20.16	4.02	667.047	645
2	100000	62.36	7.92	302.083	288.88	63.05	7.91	539.641	514.087
3	100	18.94	1.12	74.7421	71.4397	20.20	1.18	95.224	88.142
3	200	72.50	5.00	69.2354	65.8478	70.19	5.06	127.857	116.933
3	300	210.53	12.43	69.2994	65.8143	164.64	12.88	109.287	100.475
4	10000	6.16	0.96	961.75	910.571	6.76	0.96	1640.19	1524.98
4	50000	26.79	4.88	903.773	846.954	35.94	4.90	1631.8	1520.71
4	100000	77.24	10.55	824.269	776.036	77.94	10.54	1710.28	1596.1

Table 5: Computational results for A-criterion with  $K = I$ 

$\chi_i$	$n$	cpu			obj		
		mul	IP	SDPT3	mul	IP	SDPT3
1	10000	13.65	0.81	2.35	54286.3	53848.3	53848.4
1	50000	65.59	4.25	12.39	54245.2	53807.3	54103.8
1	100000	142.17	6.84	27.21	54240.1	53802.1	54103.8
2	10000	16.40	0.79	1.80	73.4521	72.4443	72.4443
2	50000	78.93	3.98	10.40	73.391	72.385	72.3853
2	100000	169.00	8.29	19.71	73.3837	72.3778	72.3777
3	100	1.48	0.95	2.13	21.6203	21.6191	21.6191
3	200	13.03	4.31	11.08	21.2826	21.2812	21.2812
3	300	37.53	9.84	30.31	21.1721	21.1706	21.1706
4	10000	12.66	0.96	1.55	174.279	170.775	170.775
4	50000	61.20	4.82	8.97	174.276	170.775	170.775
4	100000	133.68	9.48	16.43	174.277	170.775	170.776

can be efficiently computed when the method is applied to solve problem (1) with classical opti-

Table 6: Computational results for D-criterion with  $K = I$ 

$\chi_i$	$n$	cpu			obj		
		mul	IP	SDPT3	mul	IP	SDPT3
1	10000	0.92	1.05	0.86	20.5125	20.5119	20.5125
1	50000	4.52	5.19	3.81	20.5098	20.5091	20.5091
1	100000	14.89	9.77	7.56	20.5094	20.5087	20.5088
2	10000	1.78	0.81	1.07	0.410745	0.410221	0.41022
2	50000	16.99	4.37	5.06	0.409964	0.409267	0.40926
2	100000	57.42	8.65	9.05	0.409795	0.409154	0.409145
3	100	1.64	0.95	1.17	5.14292	5.14267	5.14267
3	200	16.11	4.19	6.26	5.08236	5.08212	5.08211
3	300	47.57	8.78	16.43	5.06226	5.06202	5.06201
4	10000	1.12	0.95	0.93	7.25257	7.25189	7.25189
4	50000	10.98	5.13	4.14	7.25253	7.2519	7.25189
4	100000	35.55	10.46	8.18	7.25246	7.2519	7.25189

Table 7: Computational results for  $p$ th mean criterion with  $K = I$  for some  $p \in (-1, 0)$ 

$\chi_i$	$n$	$p = -0.25$				$p = -0.75$			
		cpu		obj		cpu		obj	
		mul	IP	mul	IP	mul	IP	mul	IP
1	10000	7.63	0.99	23.3728	23.372	3.43	0.83	3635.71	3635.29
1	50000	43.86	4.47	23.3683	23.3675	25.34	4.00	3633.58	3633.2
1	100000	96.14	8.47	23.3677	23.367	60.04	7.92	3633.31	3632.94
2	10000	3.40	0.76	5.58855	5.58838	2.50	0.83	27.4836	27.4811
2	50000	21.13	3.73	5.58796	5.58771	11.49	4.08	27.4691	27.4653
2	100000	71.89	7.03	5.58785	5.58763	38.38	7.95	27.467	27.4634
3	100	1.61	0.95	6.70457	6.70448	1.53	0.99	14.1435	14.1429
3	200	14.36	3.67	6.68235	6.68225	13.45	4.08	13.9841	13.9834
3	300	42.54	8.31	6.675	6.67491	38.85	9.11	13.9318	13.9311
4	10000	1.60	0.93	7.25984	7.25955	1.30	0.96	52.2922	52.286
4	50000	9.24	4.40	7.25988	7.25956	5.68	4.79	52.2937	52.286
4	100000	30.82	8.61	7.25983	7.25957	20.84	9.03	52.2927	52.2861

Table 8: Computational results for  $p$ th mean criterion with  $K = I$  for some  $p < -1$ 

$\chi_i$	$n$	$p = -1.1$				$p = -1.2$			
		cpu		obj		cpu		obj	
		mul	IP	mul	IP	mul	IP	mul	IP
1	10000	3.92	0.75	162818	159210	3.94	0.67	485415	471459
1	50000	17.17	3.42	162740	159077	17.24	3.37	482380	471030
1	100000	47.09	7.55	162732	159060	48.57	7.21	485149	470975
2	10000	3.85	0.83	108.922	108.171	3.84	0.83	165.133	162.297
2	50000	17.98	4.06	109.588	108.072	17.05	4.08	164.314	162.133
2	100000	47.35	7.95	109.495	108.06	46.04	7.37	165.458	162.116
3	100	37.11	0.95	25.9565	25.7793	36.67	0.99	31.8264	30.8276
3	200	144.83	4.23	25.599	25.3307	143.49	4.33	31.5254	30.2362
3	300	336.23	9.56	25.5115	25.1841	334.23	9.54	31.46	30.0431
4	10000	5.85	0.98	297.604	277.597	7.39	0.92	497.138	453
4	50000	25.46	4.64	297.686	277.597	34.57	4.76	497.287	453
4	100000	65.80	9.13	297.696	277.597	81.54	9.23	497.306	453

mality criteria. Our computational results show that the IP method outperforms the widely used

multiplicative algorithm in both speed and solution quality. The codes for this paper, including our implementation of the multiplicative algorithm and our codes generating inputs for SDPT3, are available online at [www.math.sfu.ca/~zhaosong](http://www.math.sfu.ca/~zhaosong).

We would like to remark that the performance of our IP method depends on whether the Newton direction can be computed accurately and efficiently. In our implementation, we observe that for  $p$ th mean criterion with large  $|p|$ , as well as for the design space  $\{x_i = (1, s_i, s_i^2, s_i^3, s_i^4)^T, 1 \leq i \leq n\}$  with  $n \geq 50000$  and some random  $K \in \mathbb{R}^{m \times 3}$ , the Newton direction cannot be computed accurately due to numerical errors and hence our IP method fails to terminate with a good approximate solution, compared with the multiplicative algorithm.

Finally, we also compare the performance of our IP method with KNITRO [10] (Version 7.0) that is a general IP solver for solving smooth convex and nonconvex optimization problems. In particular, we call KNITRO from its Matlab interface with algorithm option `interior/CG2` to solve problem (1) for  $p$ th mean criterion with  $K = I$ . We also provide this solver with the subroutines that efficiently evaluate function, gradient and Hessian-vector multiplication. Furthermore, we use the default tolerance of the solver. We observe that our IP method is at least 10 times faster and produces a smaller objective value for all instances described in Tables 7 and 8.

## References

- [1] E. D. Andersen, C. Roos, T. Terlaky, T. Trafalis and J. P. Warners. The use of low-rank updates in interior-point methods. AdvOl-Report No. 2000/9, February 2000, Hamilton, Ontario, Canada.
- [2] A. Atkinson, A. Donev and R. Tobias. *Optimum Experimental Designs, with SAS*. Oxford University Press (2007).
- [3] A. C. Atkinson, K. Chaloner, A. M. Herzberg and J. Juritz. Optimum experimental designs for properties of a compartmental model. *Biometrics* 49, pp. 325–337 (1993).
- [4] C. L. Atwood. Sequences converging to D-optimal designs of experiments. *Annals of Statistics* 1, pp. 342–352 (1973).
- [5] C. L. Atwood. Convergent design sequences for sufficiently regular optimality criteria. *Annals of Statistics* 4, pp. 1124–1138 (1976).
- [6] C. L. Atwood. Convergent design sequences for sufficiently regular optimality criteria II: singular case. *Annals of Statistics* 8, pp. 894–913 (1980).
- [7] A. Ben-Tal and A. Nemirovski. *Lectures on Modern Convex Optimization: Analysis, Algorithms, and Engineering Applications*. SIAM (2001).
- [8] D. Böhning. A vertex-exchange-method in D-optimal design theory. *Metrika* 33, pp. 337–347 (1986).
- [9] S. Boyd and L. Vandenberghe. *Convex Optimization*. Cambridge University Press (2004).
- [10] R. H. Byrd, J. Nocedal and R. A. Waltz. KNITRO: an integrated package for nonlinear optimization. In G. di Pillo and M. Roma (Eds.), *Large-Scale Nonlinear Optimization*. Springer-Verlag (2006).
- [11] X. Chen, H. Qi and P. Tseng. Analysis of nonsmooth symmetric-matrix-valued functions with applications to semidefinite complementarity problems. *SIAM Journal on Optimization* 13, pp. 960–985 (2002).

---

<sup>2</sup>We do not use the algorithm option `interior/direct` because of the large size of the dense Hessian matrix (at least 10000 by 10000). It is not clear how to configure KNITRO to further exploit the low rank structure of the Hessian as in our IP method.

- [12] H. Dette, A. Pepelyshev and A. Zhigljavsky. Improving updating rules in multiplicative algorithms for computing D-optimal designs. *Computational Statistics and Data Analysis* 53, pp. 312–320 (2008).
- [13] V. V. Fedorov. *Theory of Optimal Experiments*. Academic Press, New York (1972).
- [14] V. V. Fedorov and J. Lee. Design of experiments in statistics. In H. Wolkowicz, R. Saigal and L. Vandenberghe (Eds.), *Handbook of Semidefinite Programming: Theory, Algorithms and Applications*. Kluwer Academic Publishers (2000).
- [15] J. Fellman. On the allocation of linear observations (Thesis). *Commentationes Physico-Mathematicae* 44, pp. 27–78 (1974).
- [16] M. C. Ferris and T. S. Munson. Interior-point methods for massive support vector machines. *SIAM Journal on Optimization* 13, pp. 783–804 (2003).
- [17] A. Forsgren, P. E. Gill and M. H. Wright. Interior methods for nonlinear optimization. *SIAM Review* 44, pp. 525–597 (2002).
- [18] R. Harman and L. Pronzato. Improvements on removing nonoptimal support points in D-optimum design algorithms. *Statistics and Probability Letters* 77, pp. 90–94 (2007).
- [19] R. Harman and M. Trnovská. Approximate D-optimal designs of experiments on the convex hull of a finite set of information matrices. *Mathematica Slovaca* 59, pp. 693–704 (2009).
- [20] R. A. Horn and C. R. Johnson. *Topics in Matrix Analysis*. Cambridge University Press (2008).
- [21] R. M. Larsen. PROPACK - Software for large and sparse SVD calculations. Available at <http://sun.stanford.edu/~rmunk/PROPACK/>.
- [22] Z. Lu and Y. Zhang. An augmented Lagrangian approach for sparse principal component analysis. To appear in *Mathematical Programming*. DOI: 10.1007/s10107-011-0452-4.
- [23] D. Papp. Optimal designs for rational function regression. Submitted (2010).
- [24] A. Pázman. *Foundations of Optimum Experimental Design*. Reidel, Dordrecht (1986).
- [25] F. Pukelsheim. *Optimal Design of Experiments*. John Wiley and Sons Inc., New York (1993).
- [26] R. T. Rockafellar. *Convex Analysis*. Princeton University Press, Princeton (1970).
- [27] S. D. Silvey, D. M. Titterton and B. Torsney. An algorithm for optimal designs on a finite design space. *Communications in Statistics – Theory and Methods* 14, pp. 1379–1389 (1978).
- [28] D. M. Titterton. Algorithms for computing D-optimal design on finite design spaces. In *Proceedings of the 1976 Conference on Information Science and Systems*, John Hopkins University, 3, pp. 213–216 (1976).
- [29] M. J. Todd, K. C. Toh and R. H. Tütüncü. On the Nesterov-Todd direction in semidefinite programming. *SIAM Journal on Optimization* 8, pp. 769–796 (1998).
- [30] K. C. Toh, M. J. Todd and R. H. Tütüncü. SDPT3 — a Matlab software package for semidefinite programming. *Optimization Methods and Software* 11, pp. 545–581 (1999).
- [31] B. Torsney. A moment inequality and monotonicity of an algorithm. In Kortanek, K.O. and Fiacco, A.V. (Eds.), *Proceedings of the International Symposium on Semi-Infinite Programming and Applications*, Lecture Notes in Economics and Mathematical Systems 215. University of Texas at Austin, pp. 249–260 (1983).



- [32] B. Torsney. W-iterations and ripples therefrom. In Pronzato, L., Zhigljavsky, A. (Eds.), *Optimal Design and Related Areas in Optimization and Statistics*. Springer-Verlag, New York, pp. 1–12 (2007).
- [33] B. Torsney and R. Martín-Martín. Multiplicative algorithms for computing optimum designs. *Journal of Statistical Planning and Inference* 139, pp. 3947–3961 (2009).
- [34] R. H. Tütüncü, K. C. Toh, and M. J. Todd. Solving semidefinite-quadratic-linear programs using SDPT3. *Mathematical Programming Series B* 95, pp. 189–217 (2003).
- [35] C. F. Wu and H. P. Wynn. The convergence of general step-length algorithms for regular optimum design criteria. *Annals of Statistics* 6, pp. 1273–1285 (1978).
- [36] H. P. Wynn. Results in the theory and construction of D-optimum experimental designs. *Journal of the Royal Statistical Society Series B* 34, pp. 133–147 (1972).
- [37] Y. Yu. Monotonic convergence of a general algorithm for computing optimal designs. *Annals of Statistics* 38, pp. 1593–1606 (2010).
- [38] Y. Yu. D-optimal designs via a cocktail algorithm. *Statistics and Computing* 21, pp. 475–481 (2011).